

LA-UR-02-1930

*Approved for public release;
distribution is unlimited.*

Title: DESIGN AND IMPLEMENTATION OF
LOW- AND MEDIUM-FIDELITY
NETWORK SIMULATIONS OF
A 30 TERAOPS SYSTEM

Author(s): Francis J. Alexander
Kathryn Berkbigher
Graham Booker
Brian Bush
Kei Davis
Adolfy Hoisie
Steve Smith

Submitted to: World-Wide Web

Los Alamos

NATIONAL LABORATORY

Los Alamos National Laboratory, an affirmative action/equal opportunity employer, is operated by the University of California for the U.S. Department of Energy under contract W-7405-ENG-36. By acceptance of this article, the publisher recognizes that the U.S. Government retains a nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or to allow others to do so, for U.S. Government purposes. Los Alamos National Laboratory requests that the publisher identify this article as work performed under the auspices of the U.S. Department of Energy. Los Alamos National Laboratory strongly supports academic freedom and a researcher's right to publish; as an institution, however, the Laboratory does not endorse the viewpoint of a publication or guarantee its technical correctness.

Design and Implementation of Low- and Medium-Fidelity Network Simulations of a 30-TeraOPS System

Francis J. Alexander, Kathryn Berkbighler, Graham Booker,
Brian Bush, Kei Davis, Adolfo Hoisie, and Steve Smith
Los Alamos National Laboratory
Los Alamos, NM 87545

Thomas P. Caudell, Donald P. Holten, Kenneth L. Summers, and Cheng Zhou
Albuquerque High Performance Computing Center
University of New Mexico
Albuquerque, NM 87131

1 April 2002

Abstract

The magnitude of the scientific computations targeted by the ASCI project requires as-yet unavailable computational power. To facilitate these computations ASCI plans to deploy massive computing platforms, possibly consisting of tens of thousands of processors, capable of achieving 10-100 TeraOPS. For various reasons the current approach to building a yet-larger supercomputer—connecting commercially available SMPs with a network—may be reaching practical limits. The path to better hardware design and lower development costs involves performance evaluation, analysis, and modeling of parallel applications and architectures, and in particular *predictive* capability. We outline an approach for simulating computing architectures applicable to extreme-scale systems (thousands of processors) and to advanced, novel architectural configurations. The proposed simulation environment can be used for: (i) exploration of hardware/architecture design space; (ii) exploration of algorithm/implementation space both at the application level (e.g. data distribution and communication) and the system level (e.g. scheduling, routing, and load balancing); (iii) determining how application performance will scale with the number of processors or other components; (iv) analysis of the tradeoffs between performance and cost; and, (v) testing and validating analytical models of computation and communication. Our component-based design allows for the seamless assembly of architectures from representations of workload, processor, network interface, switches, etc., with disparate resolutions, into an integrated simulation model. This accommodates different case studies that may require different levels of fidelity in various parts of a system. Our initial implementation, comprising low- and medium-fidelity models for the network and a low-fidelity model for the workload, can simulate at least 4096 computational nodes in a fat-tree network using Quadrics hardware. It supports studies of both simulation performance and scaling, and the properties of the simulated system themselves. Ongoing work allows more realistic simulation and visualization of ASCI-like workloads on very large machines.

1 Introduction

This report outlines the our goals and the progress we have made towards their realization as of the second quarter of FY'02. We have parallel efforts underway in the following areas:

- underlying simulation framework;
- workload models (direct and statistical);
- network models (low and medium fidelity); and,
- visualization.

The first section outlines our goals and our approach to meeting them. The next section presents the results of our evaluation of the DaSSF discrete-event engine used for the simulations. Next is a description of how network simulation can be connected to a directly executing application—an actual workload. The following section provides details on how large networks and simulations on them are visualized. The actual design and implementation of the low- and medium-fidelity network simulations follows. Finally, we outline our preliminary work involving the statistical characterization of workloads.

1.1 Motivation

The magnitude of the scientific computations targeted by the ASCI project requires as-yet unavailable computational power. To facilitate these computations ASCI plans to deploy massive computing platforms, possibly consisting of tens of thousands of processors, capable of achieving 10-100 teraOPS. For various reasons the current approach to building a yet-larger supercomputer—connecting commercially available SMPs with a network—may be reaching practical limits.

The path to better hardware design and lower development costs involves performance evaluation, analysis, and modeling of parallel applications and architectures, and in particular *predictive* capability. Performance studies are routinely used to select the best architecture or platform for a given application, select the best algorithm for solving a particular problem, and to study scalability with respect to problem and platform size. Evaluating and analyzing the performance is challenging primarily because of the large number of components making up such systems and the complex interactions that occur between them.

The tools of the trade in performance modeling and analysis are typically categorized as algorithmic/analytical analysis, statistical analysis, analysis with queuing theory, and simulation. Depending on the problem, one or more of these methods will be more appropriate than others. Although significant results have been obtained in recent work for an important class of applications of interest to ASCI [1, 2], analytical modeling of systems and applications of this scale is not always possible. Queuing models generally lead to very complex nonlinear equations whose solution is intractable. For systems of ASCI-proposed size and complexity *simulation* remains the predictive tool of choice, though simulation may be augmented by analytical and statistical analysis.

Three related targets for our simulation effort have been identified: simulation of ASCI-scale parallel systems using a realistic ASCI workload, simulation of ASCI-scale storage systems and I/O, and simulation of the high-performance ASCI wide-area network. All of these aspects of ASCI system design are equally important and tractable by the approach we propose. However, given the scale of the effort required, we envision a staged approach to tackling these problems.

In conjunction with the other methodologies, the proposed simulation environment could be used for

- exploration of hardware/architecture design space;
- exploration of algorithm/implementation space both at the application level (e.g. data distribution and communication) and the system level (e.g. scheduling, routing, and load balancing);
- determining how application performance will scale with the number of processors or other components;
- analysis of the tradeoffs between performance and cost;
- testing and validating analytical models of computation and communication such as **LogGP** [3] and **BSP** [4].

A canon of the field of performance evaluation is that hardware and software performance are inextricable—hardware performance is meaningful only in the context of applications—thus these capabilities are not entirely independent.

1.2 Goals

The *à la carte* project seeks to develop a simulation-based analysis tool for evaluating massively-parallel computing platforms including current and future ASCI-scale systems. Such a tool would provide a means to analyze and optimize the current systems and applications as well as influence the design and development of next-generation high-performance computers. Hence our general goal is to design and implement a simulation framework for design and analysis of extreme-scale parallel and distributed computing systems, and as an ongoing part of this process to validate the accuracy of results characterized by any particular model. An intermediate goal is to model (and validate the model of) the ASCI Q machine [5] with a realistic ASCI workload.

We take as given that it is not feasible to simulate an extreme scale machine and workload with perfect fidelity; on the other hand in certain circumstances it may be desirable to simulate some subset of such a machine with near-perfect fidelity. In any case the simulator itself should be able to exploit a machine of arbitrary size. Once defined, representations of logical *components* (e.g. processors or SMPs, network switches and interconnects, programs or workloads, etc.) should be easily assembled into differing configurations. Portability is essential. In more detail, the simulation system should

- be scalable to model systems comprising 10,000 processors or more;
- allow arbitrary sets of components to be represented with arbitrary degrees of fidelity, in terms of both structure (e.g. comprising distinct subcomponents) and timing;
- allow arbitrary (meaningful) configuration of components;
- allow description of the machine configuration to be as independent as possible of the descriptions of the components;
- be genuinely portable across platforms ranging from single-processor workstations to clusters of SMPs;
- be able to interface to other distinct applications such as direct execution simulators and visualization systems.

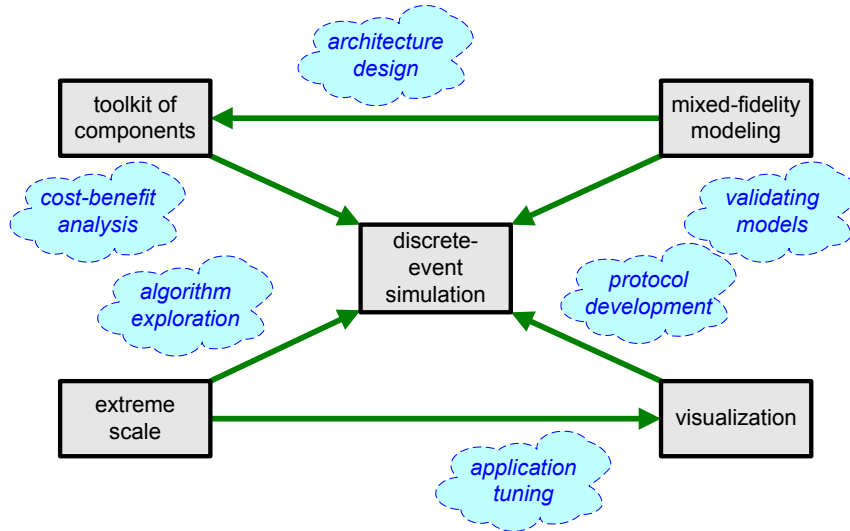


Figure 1: Goals and applications of the *à la carte* project: The boxes represent the project goals and the “clouds” represent applications for the simulation and analysis tool.

Figure 1 graphically illustrates these goals and the applications they support. These requirements suggest factoring the simulator into three parts: component descriptions, configuration descriptions, and an underlying, reasonably generic, and reasonably light-weight simulation system with which all porting issues are associated. An object-oriented approach facilitates these goals.

It is clear that simulating systems of the size and complexity that we envision will require the use of parallel simulation [6]. Furthermore, the parallel simulation substrate must support composition of simulations and be very efficient in its implementation. We concluded that a conservative synchronization scheme would have the best chance of success for this application. The requirement of portability across a variety of platforms led us to a parallel simulation substrate that runs on both shared memory and distributed memory machines. The Scalable Simulation Framework (SSF) [7] and the implementation of this framework being developed at Dartmouth College, DaSSF [8, 9], is our current choice.

1.3 Approach

Our basic approach relies on an iterative development process for constructing components of appropriate fidelities and integrating them into a portable and efficient parallel discrete-event simulation that is scalable to thousands of (simulated) computational nodes. Components may be processors, switches, network interfaces, or application workloads, for example. Studies of hardware architectures are made by running our simulation for a particular aggregate system composed of these components. The output of the simulation captures the behavior and performance of the components, and may be visualized using techniques discussed in a later section. Figure 2 illustrates the architecture of our simulator.

Simulating systems of the size and complexity we envision requires efficient parallel simulation. We use a portable, conservative synchronization engine (DaSSF), developed by Dartmouth College, for the handling of discrete events [7, 9]. DaSSF manages the synchronization, scheduling, and delivery of events in the simulation; it has a lean C++ API and supports both shared-memory and distributed-memory parallelism. We use the Domain Modeling Language (DML) to specify the

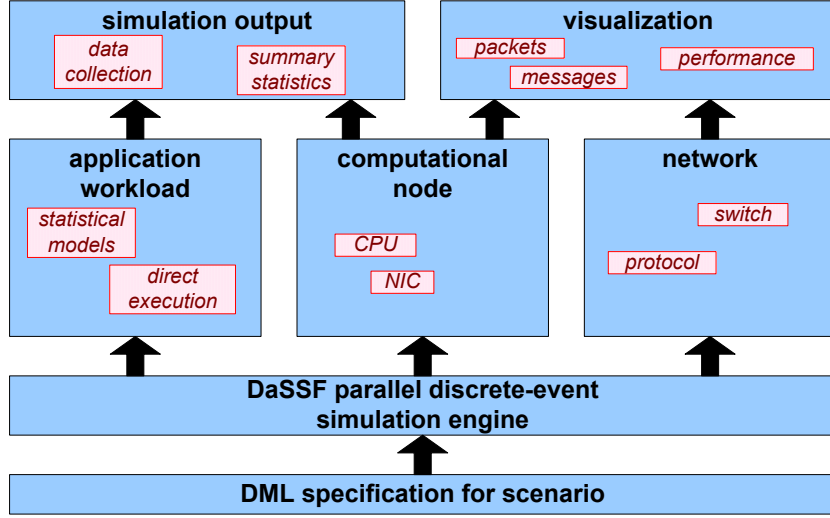


Figure 2: The architecture of the *à la carte* simulator: Simulation scenarios are represented using DML (Domain Modeling Language) and managed by the DaSSF simulation engine [9]. The application workloads, computational nodes, and networks are represented by software components that are assembled and connected according to the particular scenario being simulated. The results of the simulation may be studied visually, statistically, or in detail.

architecture and workload to be simulated. DML allows one to easily construct libraries of reusable component specifications. The lower two levels of the architecture in Figure 2 comprise DML and DaSSF.

Our component-based design (represented by the middle layer in Figure 2) allows for the seamless assembly of architectures from representations of workloads, processors, network interfaces, switches, etc., with disparate resolutions, into an integrated simulation model. One can mix and match components of different fidelities to construct a model with the appropriate level of detail for a particular study. We are focusing on the development of a simulation capability that scales to tens of thousands of processors and that can execute on a wide variety of computing platforms.

The representation of application workload (shown on the left side of the middle layer in Figure 2) forms an especially important part of the simulation. Applications and computational workloads may be represented at a variety of fidelities. Each approach below addresses tradeoffs between the accuracy of the model and the computing resources required for the model: simple *random processes* can load the hardware with message traffic having specified statistical properties. These match the distribution of messages in a real application and can include temporal and spatial correlations between messages. They ignore some of the data dependencies, however. *Direct-execution* techniques allow one to run programs nearly exactly on real processors coupled to a simulated network. These are faithful to the actual timing of an application on a processor, but may be very computationally intensive and slow. From time series of fine-grained simulations we are using learning algorithms to construct *reduced models* of the full system dynamics. This involves regression techniques like neural networks or dimension reduction methods such as the Karhunen-Loeve expansion.

The collection of simulation output (at the upper-left in Figure 2) is vitally important for understanding the behavior and performance of the simulated system. Our approach is to permit

the collection of information on all events (message sends, packetization, switching, etc.) present in the simulation at the highest level of detail. Because of the potentially voluminous nature of such data, we allow for filtering capabilities so that only data of interest will be collected in a given study. Statistical summaries also provide concise views of system performance and behavior. We are also pursuing visualization of these simulations (at the upper-right in Figure 2). We focus on both visualizing the execution of the simulation and on visualizing the performance of the simulated system.

The initial *à la carte* implementation comprises low- and medium-fidelity models of a network and a low-fidelity model of workload, but it scales to at least 4096 computational nodes in a fat-tree network. This implementation supports studies of simulation performance and scaling, and also the properties of the simulated systems themselves. Ongoing work in our iterative development approach aims to improve the fidelity of the representations and protocols. Future work will emphasize validation, the representation of I/O and storage, and wide-area networking.

2 Evaluation of DaSSF for Use in Parallel Architecture Simulation

Dartmouth SSF (DaSSF) [13, 8, 19, 9, 21] is a C++ implementation of the Scalable Simulation Framework (SSF) [8, 7]. We have completed an evaluation of DaSSF for use as the underlying parallel discrete-event-handling substrate for the *à la carte* project. Identifying a high-performance and scalable discrete-event handling mechanism is critical for the project’s success. While there may be other parallel simulation frameworks that could also meet our needs, we have focused on DaSSF and this section presents our evaluation of its suitability for our project.

2.1 Characterization of DaSSF

DaSSF is a parallel discrete event simulator that uses conservative simulation protocols to synchronize execution on multiple processors. It is the only SSF implementation that runs on both shared memory and distributed memory processors or a combination thereof. The distributed memory implementation uses MPI for communication between processors.

DaSSF is intended to be a high-performance and scalable simulator [9]. It achieves this in part by using a custom threading mechanism that handles memory very efficiently. This leads to the requirement to annotate the source code in particular ways to support the implementation of the thread behavior.

The SSF API provides five base classes that applications may subclass. The **SSF.Entity** class defines the entities in the simulation and maintains their state information. The **SSF.Process** class defines the behaviors that entities possess. Entities are connected to each other via channels, an **SSF.OutChannel** in the transmitting entity and an **SSF.InChannel** in the receiving entity. An **SSF.Event** represents the information that flows between entities across the channels. Additionally, DaSSF enhances SSF by providing classes for random number generation, data collection, and simple statistics, and for modeling semaphores, timers, and direct event scheduling.

DaSSF models may be constructed from scripts written in the Domain Modeling Language (DML). A DML script is a recursively defined list of attributes which are key-value pairs. Special keywords are included to simplify model construction by supporting the object-oriented concepts of composition and inheritance. Parsing methods are provided, but the semantic interpretation of the attributes is left to the application. The use of DML is not required; models may also be constructed programmatically from the **main** function.

DaSSF is available for a variety of platforms including the SGI Origin 2000 cluster, the Compaq AlphaServer cluster, SUN Enterprise systems, clusters of Linux workstations, and more.

2.2 Evaluation Strategy

To determine the suitability of DaSSF as a parallel simulation substrate, a small prototype model was built. The purpose was to allow us to become familiar with DaSSF and to gain experience in constructing models. The prototype was to be a learning experience and feasibility study rather than the basis for conducting a specific simulation study. The prototype that was implemented is described in more detail in Section 5.

Our requirements for the prototype were that it exercise all the essential components of DaSSF and several of the DaSSF extensions to SSF. A small model that used all the features was desirable for rapid development. At the same time we were also interested in the scaling properties of DaSSF since we want to develop very large models in the future. We were not concerned with modeling our components with high fidelity, but rather with determining whether DaSSF was an appropriate substrate for developing components with arbitrary levels of fidelity and whether these components could be easily configured into arbitrarily large models.

We were also interested in the capability of integrating DaSSF with our own C++ classes as well as with standard components such as the Standard Template Library (STL). The data collection capabilities provided by DaSSF were another topic to be investigated. The ease of debugging simulations that use DaSSF was also to be evaluated.

DaSSF is provided for several platforms and may be compiled with the native compilers or with GNU g++. We tested our model on a variety of platforms using native compilers and vendor-optimized versions of MPI.

2.3 Using DaSSF

2.3.1 Design Limitations

A concept that is important in conservative synchronization methods is that of *lookahead*. If the earliest time that a logical process at time T can schedule an event for any other logical process is at time $T + L$, then L is known as the lookahead for the first process.

A DaSSF *timeline* (another name for a logical process) is a submodel that may run concurrently with other submodels [19]. Entities are assigned to timelines, and entities on the same timeline are said to be coaligned. Each timeline maintains its own event list from which events are processed in non-decreasing time-stamp order. Entities on different timelines communicate exclusively through messages passed over channels. Entities on the same timeline also communicate via messages but may additionally use other mechanisms. In DaSSF having a substantial lookahead is more important across timelines than within a timeline. Assignment of entities to timelines is the user's responsibility, and may have a large effect on simulation performance especially when the latencies on the channels (lookaheads) differ greatly.

A simple example will illustrate how the assignment of entities to timelines affects performance. Consider four entities A, B, C and D. A is connected to B, B is connected to C, and C is connected to D. The (simulated) latencies of the connections between A and B and between C and D are much longer than the latency between B and C. In this example A and D send messages to each other with an average delay between messages of 10^6 nanoseconds and with no more messages sent after 10^7 ns. The simulation ended at 10^8 ns which was sufficient for all messages to reach their destination. Table 1 presents the results for one timeline (all entities on the same timeline), two timelines (A and B share a timeline and C and D share a timeline), three timelines (B and C share a timeline), and four timelines (each entity on its own timeline). The table includes the runtime in seconds and the number of DaSSF timeline context switches. The runs were made on a single processor Solaris workstation.

| timelines | wall secs | timeline switches |
|-----------|--------------|----------------------|
| 1 | 0.055 | 4 |
| 2 | 40.4 | 11661158 |
| 3 | 0.36 | 75030 |
| 4 | 42.3 | 12015316 |

Table 1: Example of effect of number of timelines on DaSSF performance.

| timelines | end time | wall secs | timeline switches |
|-----------|-------------|--------------|----------------------|
| 1 | 10^8 | 0.055 | 4 |
| 1 | 10^9 | 0.055 | 4 |
| 1 | 10^{10} | 0.056 | 4 |
| 3 | 10^8 | 0.36 | 75030 |
| 3 | 10^9 | 3.08 | 754455 |
| 3 | 10^{10} | 29.72 | 7347995 |
| 3 | 10^{11} | 390.88 | 75000030 |

Table 2: Example illustrating that timeline synchronization occurs even in the absence of events.

DaSSF timelines synchronize with each other using a global synchronous barrier mechanism. This synchronization occurs until the end of the simulation even when the event lists of all timelines are empty. This effect may be seen in Table 2 which reports results for the same example described above, but this time for one and three timelines and increasingly longer simulation times. Note again that all messages were completed prior to time 10^8 ns. For a single timeline the runtime remains the same regardless of when the simulation ends, while the runtime increases with increasing end time using multiple timelines.

In DaSSF the behavior of an entity is defined in one or more **SSF_Process** instances. A process that contains computations interspersed with DaSSF wait statements which are used to advance simulation time is called a procedure. DaSSF supports two types of procedures, simple procedures and procedures which are not restricted. The execution path of a simple procedure must end with a **wait** statement. Such procedures are desirable because they are implemented more efficiently in DaSSF.

The requirements for annotating the source code are not too cumbersome, but do have to be meticulously observed because DaSSF contains little error checking in this area. The DaSSF User’s Manual [9] describes the annotations, lists some limitations that derive from how the source code is parsed by the translator, and also warns of the importance of appropriately declaring **SSF STATE** variables.

DaSSF provides and manages an output collection mechanism that records output from entities distributed among any number of processors. Each processor has a single output file to which output is dumped in binary format. The **dumpData** function is a part of every **SSF_Entity**, which implies that output may only be written by existing entities. The user is permitted to define a global wrapup function which is called just before the simulation terminates. In this function one may retrieve the data from the multiple output files in time stamp order and organize and print it as desired. When the simulation terminates prematurely, such as due to an error, the data collected thus far is not post-processed by the wrapup function. This is unfortunate as the partial data could be a valuable clue to the location and cause of the abort.

The DaSSF extensions to SSF make modeling easier and we are using them despite the fact that this makes our models non-portable to other SSF implementations. Semaphores are more convenient than internal channels for signaling between processes of the same entity or coaligned entities. Timers provide a way to schedule a future action that can be cancelled at a later time if it is no longer desired. This is helpful because an **SSF.Event** cannot be canceled once it has been placed on the event list. DaSSF normally provides a process-oriented simulation world view, but it has functions that provide a discrete-event world view for situations where the extra performance gains outweigh the convenience of process-oriented simulation. We have successfully used semaphores and timers in our models. We are also using the DaSSF random number generation capabilities.

Most of the DaSSF constructs are implemented for distributed memory as well as shared memory architectures. However, two recently added features, user barrier synchronization and appointment channels, are only available for shared memory. Processes that coordinate through a semaphore must belong to entities that are coaligned in the same timeline, so these entities cannot be on distributed processors.

2.3.2 Model Building

DaSSF models may be constructed from DML files or programmatically. When DML scripts are used, the model topology is defined in the model DML file. Properties of model components, the number to be instantiated, and their connectivity are specified in this file. The machine DML file describes the hardware platform that the simulation will run on. The runtime DML file contains runtime information such as simulation start and end times and the names of the other DML files. DaSSF provides a partitioner that constructs the simulation components from the topology in the model DML file and assigns these components to the parallel computing platform.

The use of DML is not required. Models may be constructed programmatically from the **main** function. Thus far we have not experimented with this means for constructing models.

We have begun investigating the possible use of hardware description languages to formally specify architectural configurations. This will facilitate the quick assembly of architectures for case studies where an appropriate DML input file might be cumbersome to construct.

2.3.3 Portability

We regularly build DaSSF and run small simulations on a single processor Solaris workstation. We usually use the g++ compiler, but the SUNpro compiler also works fine. Any problems encountered on this most basic platform often also occur on the more sophisticated platforms, too. The only lasting difficulty we have with this platform is mentioned below where we discuss debugging issues.

Building DaSSF and running simulations on our SGI Origin 2000/IRIX cluster is often more complex than on other platforms. This occurs in part because, unlike at Dartmouth, we must interact with the cluster via the Load Sharing Facility (LSF). We prefer to use the native MIPSpro compiler and MPI implementation for better performance and local support, and after numerous attempts have a combination of options that is consistent with features required by DaSSF. We have successfully run our simple statistical models on multiple processors on multiple boxes, and we have also run a short direct execution model (described later) on this platform. A problem that remains unresolved at this time is that when a model finishes it sometimes does not properly interact with LSF to terminate the job, but rather continues to run until the LSF time limit is exceeded.

We were unable to successfully compile and run our simulation on Alpha clusters running the Tru64 operating system. Numerous problems were found using the compilers (both the native com-

ilers and various versions of g++) and libraries on these machines. Problems include segmentation faults in the compiler, name space clashes between DaSSF and the standard libraries, and failure of ISO standards-compliant C++ code to compile.

Our experience on workstations and clusters running Linux was much better. We generally had no problems compiling, linking, or running on Linux-based systems.

2.3.4 Robustness

Early on we encountered several problems in trying to use DaSSF for our prototype, but the DaSSF developer was quite responsive in addressing our questions and fixing bugs that we uncovered. We were initially unable to use DaSSF and the STL together; this has been corrected. Our early model exhibited a lot of memory leaks. This was in part due to lack of clarity in the DaSSF manual regarding which objects DaSSF automatically destroys and which are the responsibility of the user. After eliminating those leaks, we believe the DML parser still contains memory leaks. We found that DaSSF works handily with sub-directories and namespaces, except that file names and event class names must be globally unique.

There have been two major releases of DaSSF within the past year and seven bug fix releases. The documentation has been kept current with the software.

2.3.5 Debugging

We are able to debug our models using gdb on Linux machines, but have not fully succeeded in using gdb with DaSSF on Solaris where the technique that one uses to display variables on Linux does not work. Debugging is generally more cumbersome because the source-to-source translation that is needed to support the threading mechanism modifies the code and variable names. Reference to the translated code rather than the original source code is often required in order to understand the debug information.

2.3.6 Scalability

We are able to simulate 4096-node clusters using computing platforms such as the SGI Origin 2000; peak memory usage was about 8 GB in our initial prototype. Table 3 shows some of the performance figures we have obtained for a 64-node cluster. Note that performance degrades as the number of computational nodes increases because the processing nodes are not performing computational work; essentially only message passing is taking place. Simulations involving the direct execution of actual applications will have a high proportion of their CPU time spent on the workload representation, and hence will exhibit more favorable scaling characteristics.

We have also investigated how the simulation scales as a function of problem size. Figure 3 illustrates the fairly linear scaling in terms of compute time and peak memory usage for simulations of fat-trees with 8, 64, 216, 512, and 1000 computational nodes. Since the simulations were all run on a single CPU, the scaling behavior does not include contributions from additional communications overhead as would be the case if multiple CPUs were used for the larger cases. Note that the “CFD” and “uniform” workloads have similar run times, but the latter uses more memory because messages tend to queue up at the network interface cards since the network bandwidth is not sufficient to handle the demand.

The memory usage measurements discussed so far in this section were made using DaSSF prior to its version 3.2. For version 3.2, memory allocators were rewritten to avoid fragmentation. Preliminary indications are that the simulation runs described above would require only about half of the memory if version 3.2 were used. We have also noticed that the DaSSF model partitioner

| Platform | Computational Nodes | Events (1/sec) | Execution Time (sec) |
|----------------------------|------------------------|-------------------|-------------------------|
| Linux, 733 MHz Pentium III | 1 | 2020.2 | 178.1 |
| | 2 | 689.9 | 521.6 |
| Linux, 500 MHz Pentium III | 4 | 494.3 | 728.0 |
| | 8 | 379.9 | 947.3 |
| Solaris, Sun SPARC Ultra 5 | 1 | 1105.1 | 325.6 |
| Irix, Origin 2000 | 1 | 1298.3 | 277.2 |
| | 4 | 1351.1 | 266.4 |
| | 16 | 630.2 | 571.1 |

Table 3: Performance of the initial prototype simulating at 64-node architecture on a variety of computing platforms.

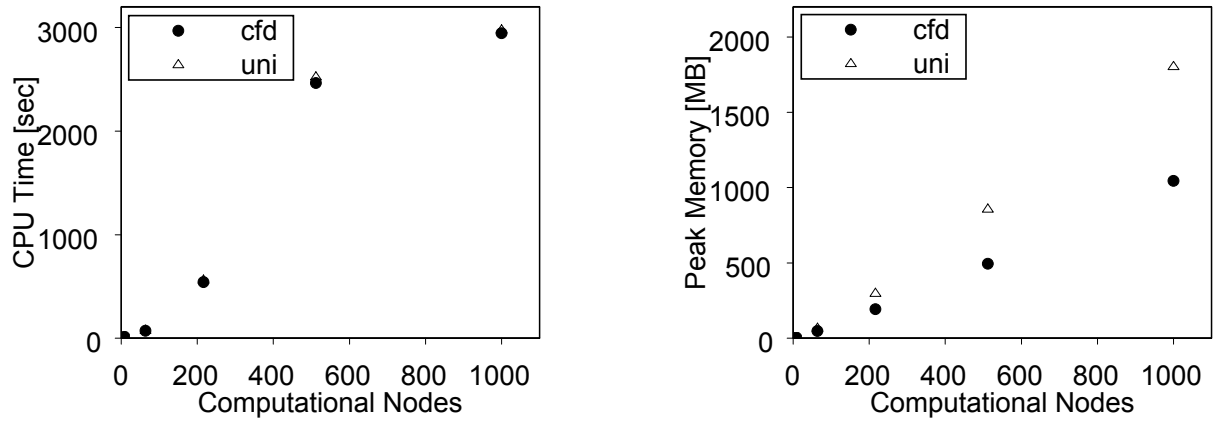


Figure 3: Scaling behavior of a fat-tree simulation on single 1 GHz Pentium III hardware running Linux: the horizontal axes are the number of computational nodes in the simulation and the vertical axes are (left) the amount of CPU time required to complete a 50 microsecond simulation with heavy message traffic and (right) the peak memory used by the simulation. The filled-in circles show results for a “CFD” workload and the hollow triangles show results for a “uniform” workload (see page 27).

tends to use disproportionately large amounts of memory if presented with a large DML input file: a DML input file containing the routing table for a 4096-node fat tree required about 10 GB of memory to partition the model.

2.4 Competing Simulation Frameworks

Although we have focused on DaSSF as our simulation substrate, we try to remain aware of other parallel discrete event simulation (PDES) environments that may be used for computer architecture simulation. A good overview of the PDES field was recently published [6]: references [27] and [20] contain surveys of languages and libraries for PDES.

SSFNet is a collection of Java SSF-based components for modeling and simulation of Internet protocols at and above the IP packet level of detail [7]. SSFNet is available in source form under the GNU public license, but it requires a Java SSF simulation kernel in order to run. A commercial implementation of Java SSF is available from Renesys Corp., and a license for research purposes is also available [7]. While we could not use SSFNet directly because it is written in Java, the common API between SSFNet and DaSSF may allow us to leverage our work with it when we reach the project stage of simulating the ASCI wide-area network.

The Wisconsin Wind Tunnel II [26], is a conservative synchronization parallel architecture simulator that utilizes direct execution to simulate multiprocessors with in-order processors and a simulated memory system. It is limited to the Sparc V8 instruction set. It appears that this simulator is no longer undergoing active development.

Researchers at UCLA have developed the Parsec language [11], and the COMPASS simulator for performance prediction of MPI programs. Parsec is a C-based parallel simulation language that supports both conservative and optimistic synchronization approaches. MPI-SIM [12], a component of COMPASS, is a library for the direct-execution driven simulation of MPI programs. It is built on top of a MPI-LITE, a library that supports multithreaded execution of MPI programs using a subset of MPI. Other components of COMPASS include a parallel I/O simulator and a parallel file system simulator [28].

A well known general purpose parallel discrete event simulation framework that uses optimistic synchronization techniques is Georgia Tech Time Warp (GTW) [10]. GTW is implemented for shared-memory multiprocessors and for heterogeneous networks of workstations, but only the shared-memory version is freely distributed [22].

Also available from Georgia Tech is Parallel/Distributed NS (PDNS) [23], a parallel implementation of the publicly available sequential *ns* simulator [25, 24]. The *ns* simulator is targeted toward networking research and includes several protocol modules. PDNS uses a conservative approach to synchronization. The current version of PDNS appears to be at least one version behind the latest release of *ns*.

3 Linking Directly Executing Workloads to the Network Simulation

The *à la carte* framework is designed to allow multiple representations for components with varying degrees of fidelity. In this section we discuss a variant of the workload component that generates interconnection network traffic according to the demands of an actual running application. This representation of the workload, known as *direct execution*, is more accurate than using simple statistical distributions to generate network traffic. (Our work on attempting to derive complex statistical characterizations of workloads, described in Section 7, falls between these two extremes.)

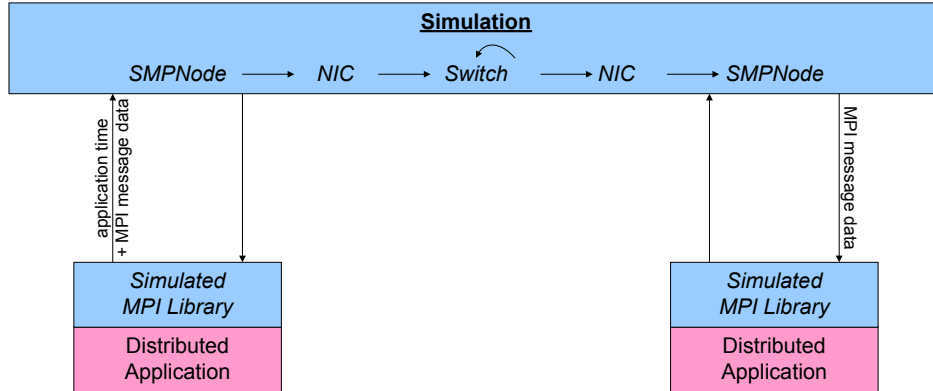


Figure 4: A distributed application directly coupled to a DaSSF-based simulation.

In direct execution simulation the application is executed on the same machine used to perform the simulation. The application is typically modified to call the simulator only for those operations that differ between the host machine and the simulated machine. Using the host machine to directly execute some instructions rather than simulating all instructions can result in considerably faster execution with minimal loss of accuracy when the host and target have similar architectures.

Our experiments with direct execution thus far have focused on simulation of communications on the interconnection network and direct execution of the computational aspects of an application. Most ASCI applications of interest use MPI (Message Passing Interface) for communication between parallel processes. We have developed a modified version of the MPI library that interfaces with a DaSSF simulation of an interconnection network. Figure 4 depicts the flow of information between the application and the simulation. Application calls to MPI are routed through the simulation rather than using the network available on the host machine. This is of course slower than using the native MPI implementation, but it allows us to study different types of interconnection networks via simulation. The modified MPI routines use timers to record the number of cycles used since the last call to an MPI routine (i.e., the amount of time spent directly executing application computations). This timing information is sent along with the MPI message data to the object that represents the sending node in the simulator. The simulator adds the application time to the simulated time before forwarding the message through the simulation to the receiving node. The final results from the application are the same as they would be without coupling to the simulation. The simulated runtime may be shorter or longer than the application runtime on the host machine depending on whether the simulated network is faster or slower than the host network.

Obtaining accurate timing information for the directly executed code is critical to this approach. We desire a timing mechanism that has high resolution to be compatible with the simulation where time is reckoned in nanoseconds, and one that is available on a variety of platforms to allow maximum portability of the simulation. Table 4 summarizes the results of our investigation into the availability of timing software packages on platforms of interest to us. The entries in the table show the resolution of the timer. Where an operating system patch is required in order to use the software, the name of the patch is indicated.

All of these timing packages have drawbacks. The fact that most counters are not saved on context switches is serious for our purposes. PAPI appears to be the most promising and is still being ported to additional platforms. The developer of PCL has indicated that it is no longer being actively developed. Other options either have limited availability or insufficient resolution.

Our experiments so far have been performed on the MIPS/IRIX platform (*nirvana.acl.lanl.gov*)

| Tool | Alpha TRU64 | Alpha Linux | Pentium Linux | MIPS IRIX | UltraSparc SOLARIS |
|---|---------------------------|----------------|-------------------------------|--------------|-----------------------|
| getrusage | μ s | ? | 10 ms | 10 ms | 10 ms |
| Atom [29] | cycle | X | X | X | X |
| PAPI [30] (patch) | cycle pfm ¹ | X | cycle perfctr | cycle | cycle cpc |
| PCL [31] (patch) | cycle pfm | X | cycle perfctr ² | cycle | cycle cpc/perfmon |
| utime [32] (patch) | X | ? | μ sec utime | X | X |
| <p>pfm: <i>pseudo device driver, not configured into kernel by default</i> perfctr: <i>Linux kernel patch for 2.2.16 to 2.4.10</i> cpc: <i>package for Solaris 8. PCL developers say it has huge overhead.</i> perfmon: <i>package for Solaris 7, need root privilege to install</i> utime: <i>Linux kernel patch for 2.2.13</i></p> <p>¹ <i>PAPI requires OSF 5.1 and pfm device driver patch. OSF 5.1 and pfm device driver does not support saving and restoring counters on context switch.</i> ² <i>PCL does not save counters on context switch.</i></p> | | | | | |

Table 4: Software packages for timing applications.

mainly because it is the only system that does not require a kernel patch to access the timing information. (Some preliminary experiments were done with Atom, but not with the application coupled to the DaSSF simulation.) We have developed a standard interface to the timing software packages that is used by our software to hide the details of accessing the individual packages, thus allowing us to easily switch between packages.

The current status of this portion of the project is as follows:

1. A replacement MPI library has been developed for those MPI calls used by SWEEP3D.
2. DaSSF processes have been developed to interface with the replacement MPI library. This involves launching the application, communicating with each distributed process through pipes, and handling each type of received message appropriately in the simulation.
3. The coupled application/simulation system has been tested on *nirvana* for a single iteration of SWEEP3D using a simple low fidelity network representation.
4. The system has been tested with both the PAPI and PCL timing packages.

Some issues need to be resolved before full scale runs with validated results can be attempted. These include:

1. Time spent in the MPI library is not presently accounted for.
2. A DaSSF problem that occurred during testing has since been addressed. The tests should be repeated with the latest version of DaSSF to verify that all components perform as expected.
3. The direct execution simulation does not terminate properly on *nirvana*. We need to understand and correct this before doing long runs there.

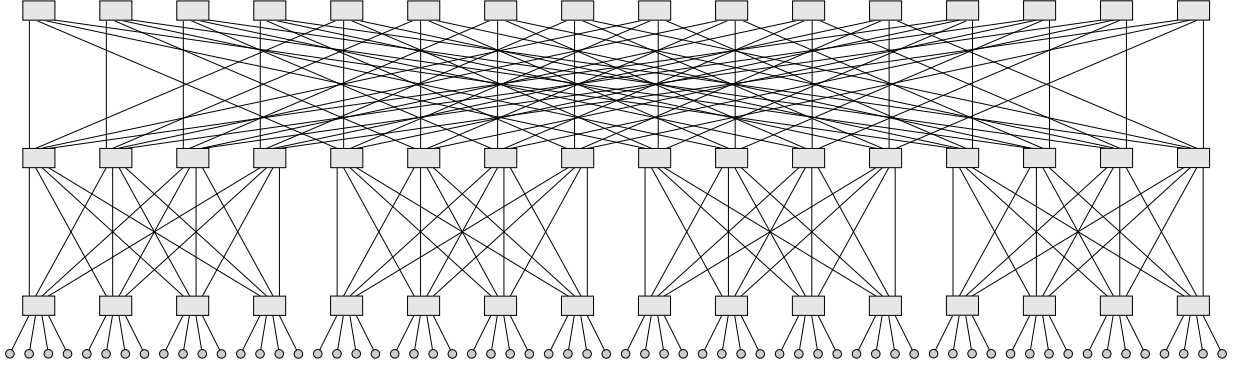


Figure 5: *Representation QS*: Quaternary fat-tree network with 64 computational nodes (small circles along the bottom) and three layers of 16 switches each (rectangles). Each switch has four duplex connections to the layer above it and four duplex connections to the layer below it—the lines in the diagram represent these connections.

4 Visualizing Simulated Fat-Tree Interconnect Networks

Our visualization efforts focus both on viewing the execution of the simulation and on displaying the performance of the simulated system. Visualization also aids in debugging the simulation itself, in developing and evaluating the efficiency of load balancing of the simulation entities, and in understanding synchronization between simulation timelines. Visualizing the simulated system allows end-users to understand how varying workload or network architecture affects the overall performance of an advanced or novel architecture. They can also see communication patterns, levels of network usages, and the presence of bottlenecks. Our visualization approaches include direct representations of the architecture as well as innovative abstractions of the architecture and dynamics of the system.

For simplicity we consider so-called “single-rail” networks where each computation node only has one network interface card (NIC). We also assume each network switch has eight duplex I/O ports; the ports may be linked to computational nodes or to other switches. The network is organized into layers of switches that connect only to the layers above and below: for eight-port switches, we have four upward connections and four downward connections. Thus we have a quaternary fat-tree network: “quaternary” because each switch has four-fold connections upwards and downwards, and “fat” because the number of switches per layer is the same for all layers. (Note that the real-life networks are typically “thinned” at the higher layers by reducing the number of switches there.) Figure 5 illustrates the layout of such a network with 64 computational nodes and three layers of 16 switches each.

We start with some formalism describing these networks. Let L be the number of layers in a complete (i.e., not thinned) quaternary fat-tree network. Number the layers $\ell = 1, \dots, L$. Each layer has $4^{L-\ell}$ switches, so we label these with IDs $x = 0, \dots, 4^{L-\ell} - 1$. This means there are $s_L = L \cdot 4^{L-1}$ switches in the network and $n_L = 4^L$ nodes connected to it. Each port p of switch x in layer ℓ connects to four switches

$$y_p = 4^\ell \left\lfloor \frac{x}{4^\ell} \right\rfloor + 4^{\ell-1}p + (x \bmod 4^{\ell-1}) \quad (1)$$

in layer $\ell' = \ell + 1$ at port $p' = p + 4$, where $p = 0, 1, 2, 3$.

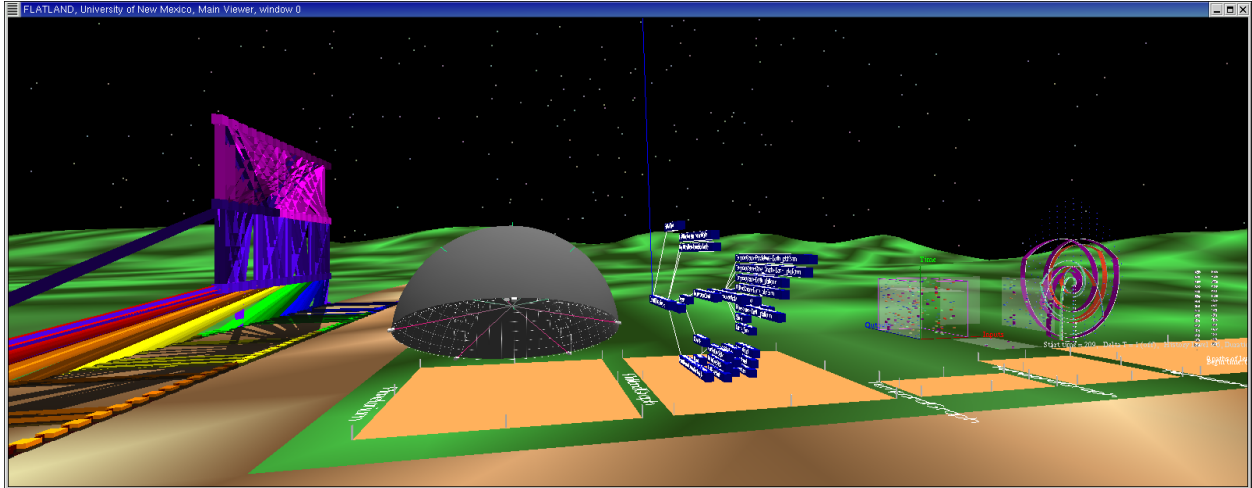


Figure 6: Flatland with a variety of independently developed applications from various problem domains.

4.1 Flatland: An Immersive Visualization Development Framework

Flatland is an immersive visualization development framework created at the University of New Mexico as part of the Homunculus project [33]. It is used to facilitate rapid prototyping and research in scientific and information visualization, immersive environments and interfaces, and human factors engineering. The Flatland infrastructure aids in: creating and managing complex scene graphs with OpenGL geometry, lighting, shadows, stereo rendering, and spatialized sound objects; dynamically loading applications in ways that allow them to interact without interfering; managing novel input and output devices; facilitating navigation of the resulting virtual spaces; and some basic, optional “world” objects such as the landscape, stars and sun that you see in Figure 6.

One of the tenets of our approach to visualization is that immersive, 3D environments can offer unique advantages over two dimensional plots, charts or graphs as well as simple, non-immersive three-dimensional graphics. By placing the user of these tools within the same context as the objects being viewed and allowing the user to navigate around them, picking their own points of view and using motion parallax to comprehend the three dimensional relationships between objects, allowing them to cast shadows and perhaps even to have behavior and emit sounds, a qualitative improvement in comprehension of the data is achieved.

4.2 Direct Representation

For purposes of debugging a simulation, it is useful to have a visual representation that clearly and explicitly represents all of the simulation’s components. With these fat-tree-connected massively parallel computing architectures, we started by laying out the processors, NICs, and switches in several relatively obvious direct representations. Figure 7 shows a three-dimensional visualization of a 64-node network where the layers are arranged vertically so that the computational nodes are on the bottom and the successive layers of switches are arranged above one another. The first (Representation DL) is a simple distribution of processors and NICs along a line with the layers of switches above them also along a line. The second (Representation DC) simply involves wrapping this line around into a circle, creating a sort of cone topped with a cylinder. The third (Representation DR) is a rectangular grid laydown of the processors, NICs and switches. The representations shown are for a 64 processor machine with 3 layers of 16 switches each. Only the

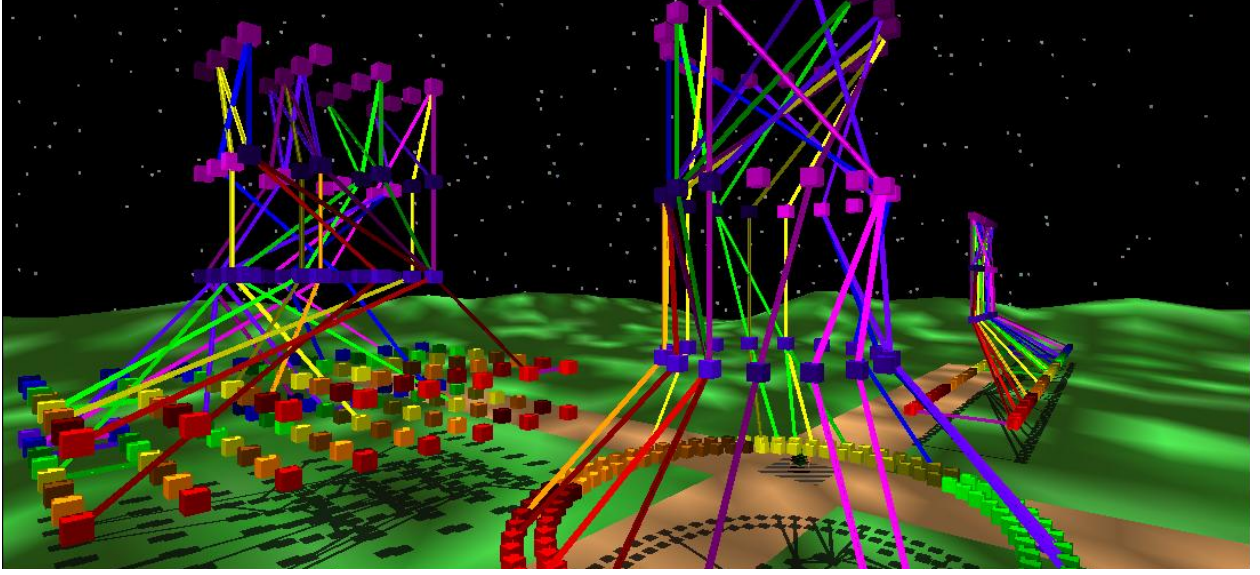


Figure 7: *Representation DR* (left), *Representation DC* (right, front), and *Representation DL* (right, rear): The bottom layer of boxes represents computational nodes and the upper layers of boxes correspond to the layers of switches. Only the links for messages currently in the network are drawn. The boxes and links are color-coded according to their attributes and the whole visualization is animated as simulation time progresses.

links for messages actively being transmitted through the network are drawn (and color coded according to message properties). The visualization is animated so that the user can see the progression of messages sent through the network.

Although the direct representations provide the highest level of detail concerning message properties and component connectivity, they suffer from the lack of scalability and the visual complexity of the display (i.e., number of overlapping lines) when large models (more than several hundred computational nodes) are displayed. The *ad hoc* attribute coloring capabilities make this representation extremely valuable for debugging purposes, however. The largest machine we have studied is a 4096 processor machine with 6 layers of 1024 switches each. These representations were overwhelmed at that scale.

4.3 Layered Block Representation

In order to address the scalability and visual complexity issues of the direct representations, we have developed a more abstract representation (motivated by the abstraction of a graph as a connection, or adjacency, matrix) where switch layers are grouped together into aggregate visual objects. The whole network is laid out on a square with equally-spaced pillars along the diagonal representing the computational nodes and their NICs. The various switch layers are grouped by fours on succeeding levels below. For example, the first level below the nodes consists of groups of four switches, the second level consists of groups of sixteen switches, etc. Figure 8 shows a view of this representation. When a processor sends a message to another processor through the switching fabric, a line or “pipe” leaves the processor and grows across the switching fabric until it comes to the point on the implied connectivity matrix where the two processors’ connection would normally be indicated. At that point it makes a 90 degree turn and continues to grow until it reaches the destination processor. In one mode of use, the layered switch blocks change color as

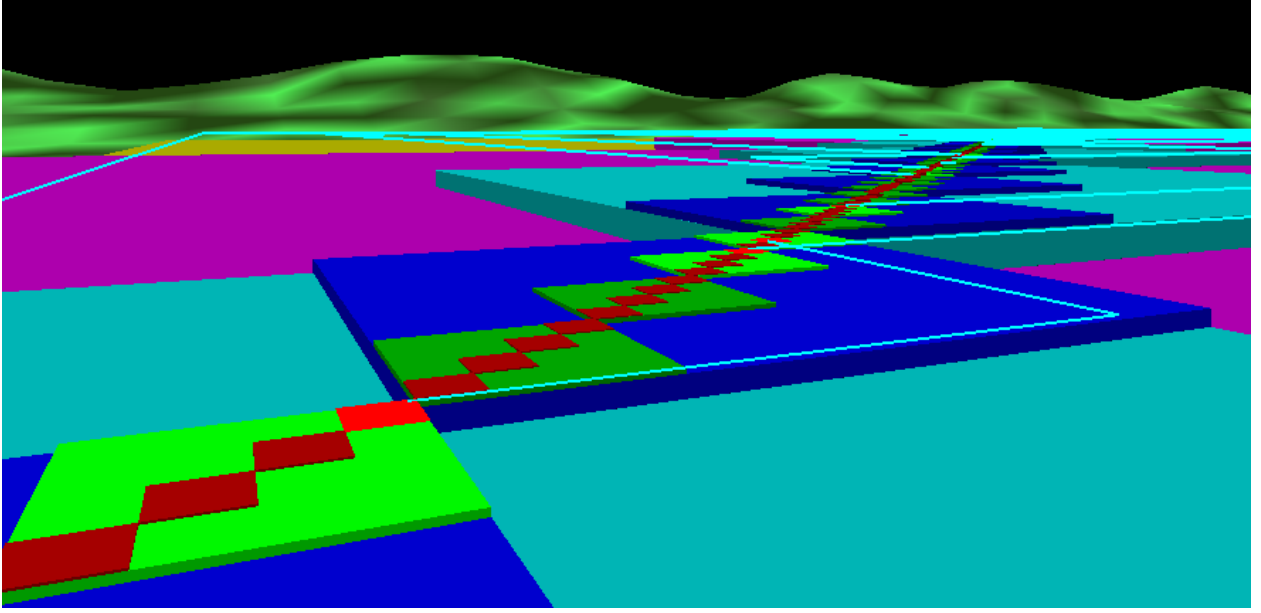


Figure 8: *Representation LB*: The network is laid out in a square. The first layer of switches is below this in groups of four (shown in red), the second layer of switches is below that in groups of sixteen (shown in green)—subsequent layers are blue (3rd), aqua (4th), magenta (5th), and yellow (6th). The switch groups brighten in color when one of their switches is active, and the “pipes” represent individual messages passing through the switches below them.

they are involved in more and more communication traffic, allowing the viewer to recognize the relative level of utilization of each switch or switch group. It is also possible to subdivide a switch group vertically into individual switches.

The layered representation provides a more compact display than the direct representation and is somewhat more scalable, but it still suffers from similar problems once the number of nodes approaches one thousand. It does have the advantage of being able to show the activity in very large systems (e.g., 4096 computational nodes) if the pipes showing individual messages are suppressed. Unfortunately, the connectivity matrix scale is of order n_L , which naturally limits its scalability. At 4096 processors, individual processors, NICs, and even the first level of switches are smaller than a single pixel when viewed in their entirety, even on a high resolution screen (1600×1200). This is only a partial limitation since two dominant modes of use are likely to be macroscopic—attempting to understand the aggregate dynamics of the system which will be mostly differentiated at higher than NICs and first level of switches—or microscopic—focusing on following single messages through the system. Nevertheless, in this representation, full detail cannot be apprehended in full context simultaneously.

4.4 General Framework for Fat-Tree Representations

We now consider a general framework that allows us to express the layout of computational nodes and network switches in a two- or three-dimensional visualization of a fat-tree. It is actually sufficient to consider only switches in a layout, as each quadruplet of nodes connects to only a single switch—the lowest-level switch in any layout can be replaced by that switch and its four connected nodes. Likewise, we need not consider the eight duplex ports on each switch because each switch could be expanded into its eight *in* ports and/or its eight *out* ports in a layout.

Because some of these layouts have a fractal nature, we make a brief diversion to discuss the measurement of fractional dimension in these representations. Now we limit our consideration to two-dimensional layouts whose aspect ratio does not vary with L . Let w_L be the width of the layout in pixels. The box-counting dimension d of the layout is defined from $s_L \sim w_L^d$ as $L \rightarrow \infty$. To compute this, we take a limit:

$$d = \lim_{L \rightarrow \infty} \frac{\log s_L}{\log w_L} = \lim_{L \rightarrow \infty} \frac{\log L + (L-1) \log 4}{\log w_L} = \lim_{L \rightarrow \infty} \frac{(L-1) \log 4}{\log w_L}. \quad (2)$$

Here is a general technique to represent the plethora of possible fat-tree layouts: Consider a pair of generating functions A_n and B_n which map the integers 0, 1, 2, 3 to pixel coordinates; here n is a non-negative integer specifying the scale of the mapping. It is very important that the range of the functions A_n and B_n are disjoint—otherwise, switches on different layers will overlap on the same pixel. Since the fat tree is quaternary, it is useful to represent switch IDs in base four: namely, represent a number x as $x = \sum_{n=1}^{L-1} 4^{n-1} x_n$ where the x_n are its base-four digits. We can now write the pixel coordinates of switch x in layer ℓ out of a total of L layers as:

$$F_{L,\ell}(x) = \begin{cases} (0, 0) & \text{for } L = 1 \\ \sum_{n=\ell}^{L-1} A_n(x_n) & \text{for } \ell = 1 \\ \sum_{n=\ell}^{L-1} A_n(x_n) + \sum_{n=1}^{\ell-1} B_n(x_n) & \text{for } 1 < \ell < L \\ \sum_{n=1}^{L-1} B_n(x_n) & \text{for } \ell = L \end{cases} \quad (3)$$

The layout is generated as the union of all the pixel coordinates of the switches:

$$\bigcup_{\ell=1}^L \bigcup_{x=0}^{4^{L-1}-1} F_{L,\ell}(x). \quad (4)$$

Note that although we have considered two-dimensional layouts here, the formalism extends trivially to the three-dimensional case.

4.5 Compact, Self-Similar, and Fractal Representations

At this point it should be clear that representations which essentially scale linearly with n_L are going to fail in at least the most obvious way of not providing full detail in full context. While level of detail management would help make this *dynamic range* problem somewhat more graceful, it would not allow us to apprehend each individual node, NIC, and switch, all at the same time without improving the *compactness* of our layout. The *direct conelike* representation described in Section 4.2 scales by approximately n_L , which is still linear with n_L but the *direct rectangular* layout scales by the $\sqrt{n_L}$. Thus 4096 processors, for example, only require an area on the order of 64×64 units to display. This seems very promising but by distributing the processors in a 64×64 array, the majority of the processors are in the middle of the area and similarly, the switching layers, laid out in 32×32 arrays tend to occlude each other nearly as badly. From this simple analysis, it appears that laying the processors, NICs, and each switch layer out in two dimensions is compact enough for our needs but leads immediately to occlusion problems.

Motivated by the somewhat self-similar nature of the fat tree, we began to investigate two different representations, one inspired by a simple pair of fractal generators as described in Section 4.5.1 and another inspired by *H-array radar antennae* as described in Section 4.5.2 below. The similarities between these two representations lead us to consider a more general representation of all layouts in two and three dimensions of fat trees.

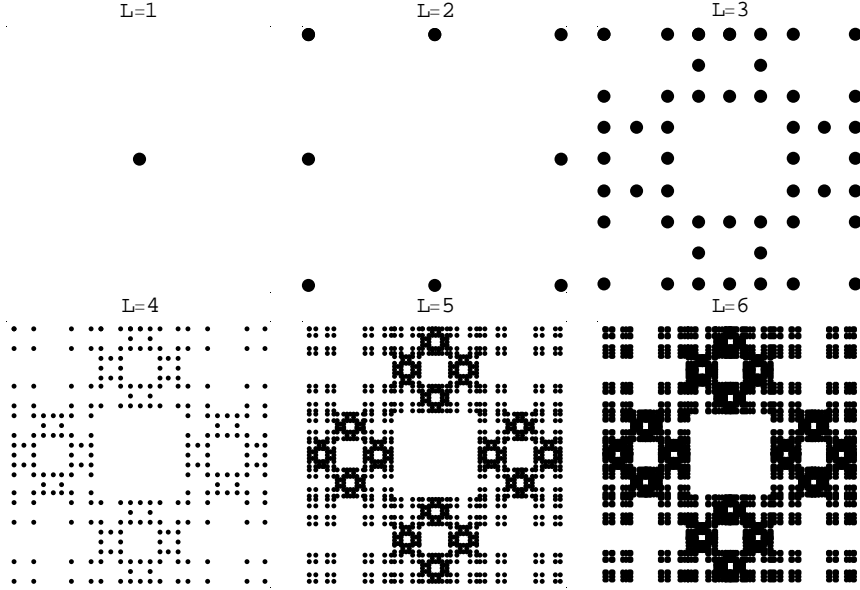


Figure 9: *Representation F1*: The six panels show placement of switches (circles) for the representation generated by Eq. (5) for $L = 1, \dots, 6$.

4.5.1 Fractal Representation

We start with a fractal-based representation defined by

$$A_n(k) = 3^{n-1}a(k) \text{ and } B_n(k) = 3^{n-1}b(k) \quad (5)$$

where

$$a(k) = \begin{cases} (0, 1) & \text{for } k = 0 \\ (1, 0) & \text{for } k = 1 \\ (0, -1) & \text{for } k = 2 \\ (-1, 0) & \text{for } k = 3 \end{cases} \text{ and } b(k) = \begin{cases} (-1, 1) & \text{for } k = 0 \\ (1, 1) & \text{for } k = 1 \\ (1, -1) & \text{for } k = 2 \\ (-1, -1) & \text{for } k = 3 \end{cases}. \quad (6)$$

The function $a(k)$ places the lower-layer switches on the sides of a square, while the function $b(k)$ places the higher-layer ones on the corners of the same square. The 3^{n-1} coefficient in Eq. (5) ensures that subsequent squares are appropriately scaled to a larger size. Figure 9 shows the fractal for the first several L —one can plainly see the self-similarity between networks of different sizes.

It is easy to see that $w_1 = 1$ and that $w_L = 3w_{L-1}$ for this representation. The solution to this recursion relation is $w_L = 3^{L-1}$, which results in a fractal dimension $d = \log 4 / \log 3 \approx 1.26$.

This layout has the advantage that it scales well—one can represent the 6144 switches of a six-layer fat-tree in a 243×243 pixel area, for instance. It has a disadvantage that the switches for different layers are interleaved, making it somewhat tricky to visually separate the activity in different network layers. We have used animated versions of these layouts to successfully distinguish the distribution of messages in two 4096-computational-node applications with different communication patterns, however. There are several variations of simple animation that might be useful with this representation:

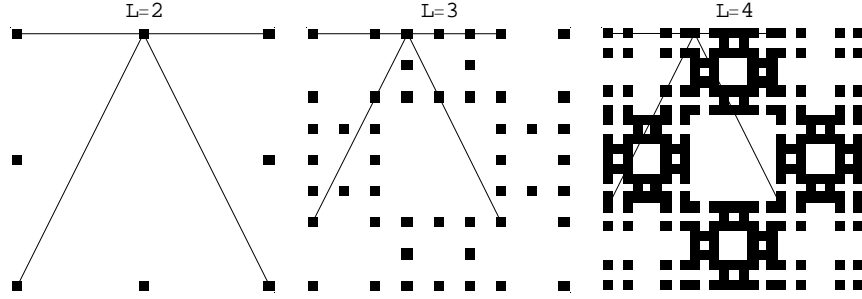


Figure 10: Connection patterns between consecutive layers of switches in Representation F1 (Fig. 9) for $L = 2, \dots, 4$.

- Color the switches by the ID of the message they are transmitting.
- Color the switches by the number of ports they have in use.
- Add a time axis as the third dimension to convert the animation into a static volume visualization.
- Add a message ID axis as the third dimension to convert the animate into a dynamic volume visualization.
- Expand each switch so that each of its eight ports can be seen.
- Display computational nodes in addition to switches.
- Use the third dimension to animate connections between switches with arch-like pipes.

We can also examine the patterns of connections between switches at consecutive layers in this layout. Figure 10 illustrates this self-similar linking of switches.

4.5.2 “Fat H” Representation

We can address the problem of switch layers being interleaved in Representation F1 in a new representation defined by

$$A_n(k) = (n + 2)2^{n-2}b(k) \text{ and } B_n(k) = 2^{n-2}b(k) \quad (7)$$

where

$$b(k) = \begin{cases} (-1, 1) & \text{for } k = 0 \\ (1, 1) & \text{for } k = 1 \\ (1, -1) & \text{for } k = 2 \\ (-1, -1) & \text{for } k = 3 \end{cases} \quad (8)$$

as in the other representation. In this case both the lower and upper level switches lie on the corners of a square, but the coefficient $(n + 2)$ in Eq. (7) forces the lower level switches onto the corners of a larger square. Thus the more central groups of switches are in higher layers. Figure 11 illustrates this.

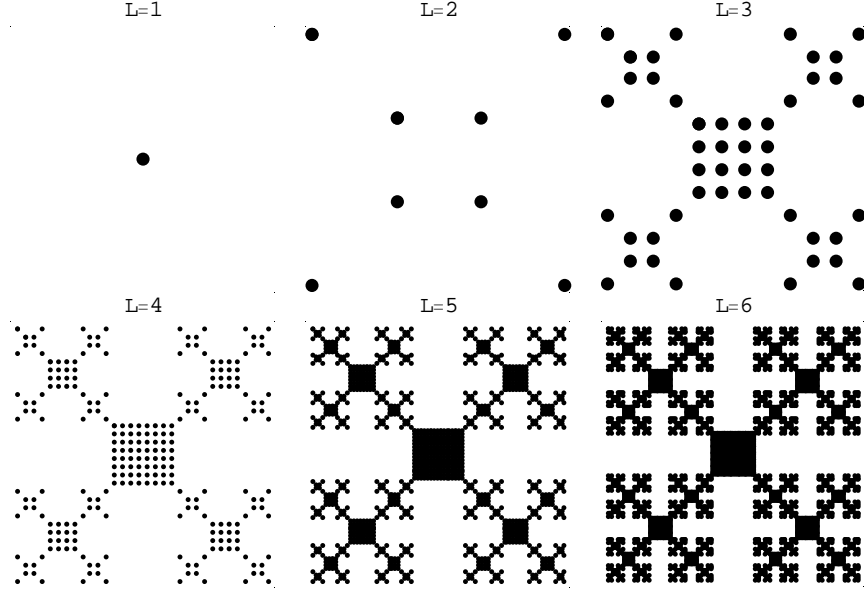


Figure 11: *Representation HT*: The six panels show placement of switches (circles) for the representation generated by Eq. (7) for $L = 1, \dots, 6$.

It is easy to see that $w_L = 1$ and that $w_L = 2w_{L-1} + 2^{L-1}$ for this representation. The solution to this recursion relation is $w_L = L \cdot 2^{L-1}$. This results in a fractal dimension $d = 2$. Hence, this representation “efficiently” fills two-dimensional space and is not technically a fractal (since it does not have fractional dimension). One can see this in Figure 11 in that the highest, central layer of switches occupies a smaller region of the diagram relative to the layers below as L increases. Hence, the layout is not self-similar as a function of L . One could rectify this situation by altering the scaling factors in Eq. (7) at the expense of letting higher layers take relatively more area in the diagram.

The connection pattern between switches possesses a more straightforward arrangement in Representation HT (Fig. 12) than for the more interleaved Representation F1 (Fig. 10).

4.5.3 Additional Representations

We can also reformulate other representations in terms of the formalism of Equation (3). Representation QS (Fig. 5) can be expressed as

$$A_n(k) = \left(\left(k - \frac{3}{2} \right) 4^{n-1}, -\frac{1}{2} \right) \text{ and } B_n(k) = \left(\left(k - \frac{3}{2} \right) 4^{n-1}, \frac{1}{2} \right). \quad (9)$$

Representation LB (Fig. 8) can be written, more or less, as

$$A_n(k) = 4^{n-1} \left(k - \frac{3}{2}, k - \frac{3}{2} \right) \text{ and } B_n(k) = 4^{n-1} \left(k - \frac{3}{2}, (k + 2 \bmod 4) - \frac{3}{2} \right). \quad (10)$$

In general, to create additional representations, one only needs to choose a pair $(A_n(k), B_n(k))$ whose ranges do not overlap.

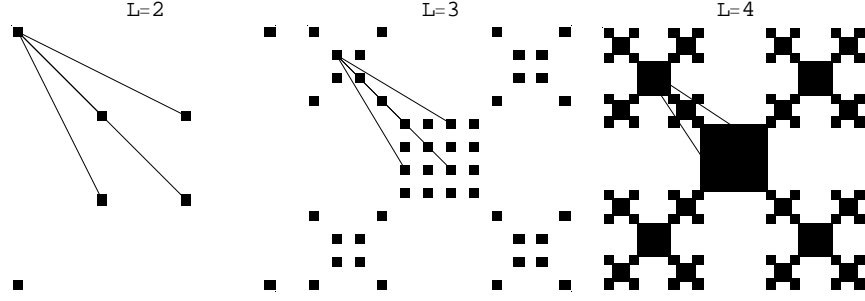


Figure 12: Connection patterns between consecutive layers of switches in Representation HT (Fig. 11) for $L = 2, \dots, 4$.

| L | s_L | n_L | $w_L^{(F1)}$ | $s_L / \left(w_L^{(F1)}\right)^2$ | $w_L^{(HT)}$ | $s_L / \left(w_L^{(HT)}\right)^2$ |
|-----|-------|-------|--------------|-----------------------------------|--------------|-----------------------------------|
| 1 | 1 | 4 | 1 | 1.00 | 1 | 1.00 |
| 2 | 8 | 16 | 3 | 0.89 | 4 | 0.50 |
| 3 | 48 | 64 | 9 | 0.59 | 12 | 0.33 |
| 4 | 256 | 256 | 27 | 0.35 | 32 | 0.25 |
| 5 | 1280 | 1024 | 81 | 0.20 | 80 | 0.20 |
| 6 | 6144 | 4096 | 243 | 0.10 | 192 | 0.17 |

Table 5: Comparison of sizes of the layouts for Representations F1 and HT.

4.5.4 Comparison

It is instructive to compare the advantages and disadvantages of some of the fat-tree representations we have created so far. Table 5 contrasts the dimensions of Representations F1 and HT. Thus these have comparable sizes for the networks of interest ($L \leq 6$). Table 6 ranks all six representations in terms of usability metrics.

5 Low-Fidelity Quadrics Network Simulation Design and Implementation

The initial *à la carte* low-fidelity implementation consists of compute nodes with a simple statistical model of workload (exponentially-distributed message sizes and spacings) coupled to a circuit-switched fat-tree network. The source-target patterns for messages can be configured. The switches with four *up* ports and four *down* ports are modeled at the packet level with a simple protocol that

| Property | QS | DL | DC | DR | LB | F1 | HT |
|--------------------|-----------|------|----------|------|----------|----------|----------|
| Microscopic Detail | very | very | very | very | yes | somewhat | somewhat |
| Scalability | no | no | somewhat | yes | somewhat | yes | very |
| Visual Complexity | very high | high | moderate | high | moderate | moderate | low |

Table 6: Advantages and disadvantages of various representations.

allocates a circuit through the network for each message. We have run the simulation on a variety of sample models from 1 to 4096 computational nodes (up to 6144 switches).

5.1 Requirements

Our requirements for the low-fidelity implementation were that it exercise all the essential components of DaSSF and several of the DaSSF extensions to SSF. To investigate the scalability of DaSSF itself the components of the implementation were to be such that they could be easily configured into arbitrarily large models. Since we were not conducting a real performance study, we were not concerned with modeling our system components with high fidelity, but rather determining whether DaSSF was an appropriate substrate for developing components with arbitrary levels of fidelity. A small model that used all the features was desirable for rapid development. At the same time we were also interested in the scaling properties of DaSSF since we will want to develop very large models in the future. The data collection capabilities provided by DaSSF were another topic to be investigated.

DaSSF is provided for several platforms and may be compiled with the native compilers or with g++. We tested our model on a variety of platforms using native compilers and vendor-optimized versions of MPI. We were also interested in the capability of integrating DaSSF with our own C++ classes as well as with standard components such as the Standard Template Library. The ease of debugging simulations that use DaSSF was also to be evaluated.

5.2 Design

We used the results of our domain analysis to define a subset of components to implement in the implementation. We chose not to model the processors and memory hierarchy of an SMP node in any detail and instead to focus on modeling the interconnection network between nodes as a fat-tree network with a circuit-switched routing protocol. For the workload, rather than model the message traffic from any specific application, we chose to have each processor node emit messages with exponentially distributed interarrival times. The message destination node is selected with uniform probability, and the message size is exponentially distributed. The parameters for the distributions are inputs to the simulation.

The DaSSF API provides five base classes that applications may subclass. The **SSF_Entity** class defines the entities in the simulation. The **SSF_Process** class defines the behaviors that entities possess. Entities are connected to each other via channels, an **SSF_OutChannel** in the transmitting entity and an **SSF_InChannel** in the receiving entity. An **SSF_Event** represents the information that flows between entities across the channels. Parameterized properties of model components, the number to be instantiated, and their connectivity are specified via an input file written in the Domain Modeling Language (DML). A simplified UML (Unified Modeling Language) diagram of our low-fidelity model is shown in Figure 13. The model contains three types of entities, representing the SMP node, the network interface card (NIC), and the network switch.

The SMP node has an outgoing channel for sending messages to its NIC and an incoming channel for receiving messages from its NIC. The NIC has a corresponding incoming channel for receiving messages from its SMP node and an outgoing channel for sending messages to its SMP node. Additionally, the NIC has an outgoing channel and an incoming channel that connect it to its network switch. Each network switch has 8 incoming channels and 8 outgoing channels. At the level nearest the NICs, 4 of the channels communicate with 4 NICs and 4 communicate with the next level in the fat tree. (An example of a fat-tree network is shown in Figure 2). Higher in the fat tree communication involves only other switches. The route that a message takes through the

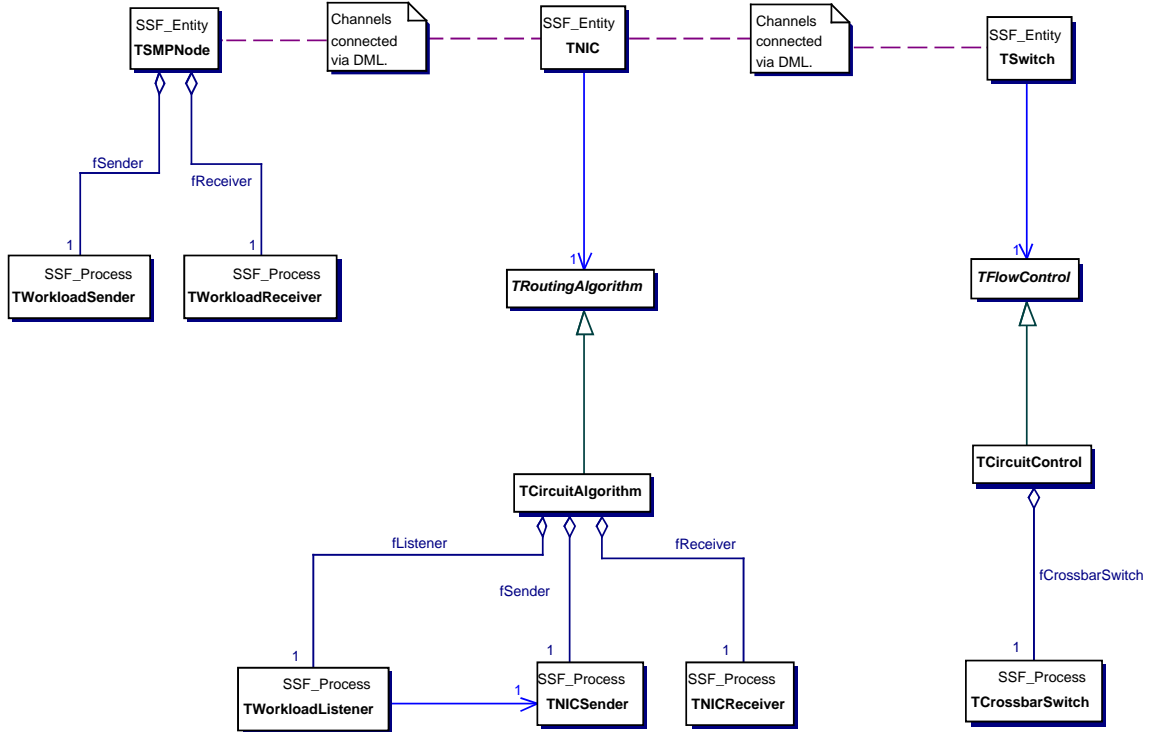


Figure 13: Simplified UML (Unified Modeling Language) class diagram for the low-fidelity simulation.

network is determined at the source and is the same for all packets in a message. The route may be read and stored in a routing table or computed via a routing method.

The SMP node has two processes, **TWorkloadSender** for sending messages and **TWorkloadReceiver** for receiving messages. The NIC has a routing algorithm which has three processes, **TWorkloadListener** for receiving messages from the source SMP and buffering them if the NIC is busy, **TNICSender** for splitting the message into packets and sending the packets to the network switch, and **TNICReceiver** for receiving packets from the switch, reassembling the message and sending it to the destination SMP. The network switch has a flow control algorithm which has one process, **TCrossbarSwitch**, that receives packets on an incoming channel from a NIC or another switch and forwards them out the appropriate outgoing channel to the next switch or NIC as specified by the route that is embedded in the packet.

We defined separate processes for the actions performed by the SMP node and NIC entities as a way of modularizing the logic, to reflect the fact that a NIC may represent a significant process running asynchronously with the SMP, and to allow for the possibility of multiple NICs per SMP. However this necessitates the use of other mechanisms such as a semaphore or an internal channel to allow processes belonging the same entity to communicate with each other. Rather than attaching the processes for message and packet handling directly to the NIC, we interposed the abstract **TRoutingAlgorithm** object and created a **TCircuitAlgorithm** class that contains the processes which implement a circuit routing algorithm for packet delivery. This gives us the flexibility to define other algorithms in the future and give new behavior to the NIC by simply selecting at runtime a different algorithm to be constructed. Similar logic pervades the design of the switch and its process.

Two types of events are defined, one for a message, and one for the packets in a message. Packets

are further subclassed into four types: the open circuit packet, the acknowledge circuit packet, the close circuit packet, and the data packet. The precise behavior that occurs in the processes depends upon the type of packet that arrives.

Our design employs the DaSSF-supplied random number generation module. We collect traces of packet movement through the network using DaSSF's data packing and dumping capability.

5.3 Implementation

The present implementation of our low-fidelity simulator consists of approximately five thousand lines of ANSI- and POSIX-compliant C++ that runs using MPI under the Linux, Solaris, and Irix operating systems. We use DaSSF's built-in user threads to avoid the operating-system overhead associated with creating hundreds of native threads. We have executed the code on single processor boxes, SMPs, and clusters of workstations on a LAN. Since the specifications of modelled components and architectures are entirely independent of the host platform, models may be developed and executed on any convenient or appropriate host.

The data-collection facility in the implementation permits one to collect detailed information concerning the history of each simulated message: its movement from the computational node to the network interface card, the opening of a network circuit for its transmission, its packetization, its acknowledgement, and the closing of its network circuit. The implementation has been tested to verify correctness of time delays, probabilistic distributions, and synchronization behavior. The data-collection also supports our visualization capability, allowing users to view message traffic in the simulated system. Data summaries supply information on queue sizes, throughputs, port usage, timeouts, path lengths, and communication patterns.

5.4 Results

We have constructed a variety of architecture models for testing the simulation and for measuring its performance and scalability. The simplest models consist of computational nodes connected to each other by a bus or via a pair of network interface cards; slightly more complex models contain one or two network switches with eight computational nodes. These small models provide a testbed for verifying that the timing delays in various parts of the simulation match those of the hardware being simulated. They also supply a convenient platform for debugging and detailed tracing of the simulation's progress. Even a simulation of a small system executing for a short time can produce a large amount of data if the history of each simulated entity, process, and event is recorded.

Our larger models contain 64, 128, . . . , up to 4096 computational nodes networked together in a fat-tree topology. Figure 5 shows the layout of a 64-node system. A 4096-node system (similar in size to a proposed ASCI Q machine architecture) uses these 64-node systems as building blocks—the nodes of one such system are each replaced with a 64-node system. Having these more realistically-sized models allows us to undertake meaningful studies of system performance and scaling—both of the simulation system itself and of the architectures being simulated. Our current studies focus on the behavior of the simulation system rather than the simulated architecture. In Section 2.3.6, above, we discussed the scaling properties of this simulation.

Figure 14 shows the results of a simulation of the 64-node, 48-switch architecture shown in Figure 5, with a load of 100 MB/s of message traffic originating at each node. The diagram shows that many of the messages reach their destination in the minimum possible time (based on network connectivity and time delays), but that a significant fraction of the messages have additional delays due to network congestion. The mean message size is 4 KB. The simulation is not realistic in the sense that the network protocols do not correspond to those of any actual network hardware and

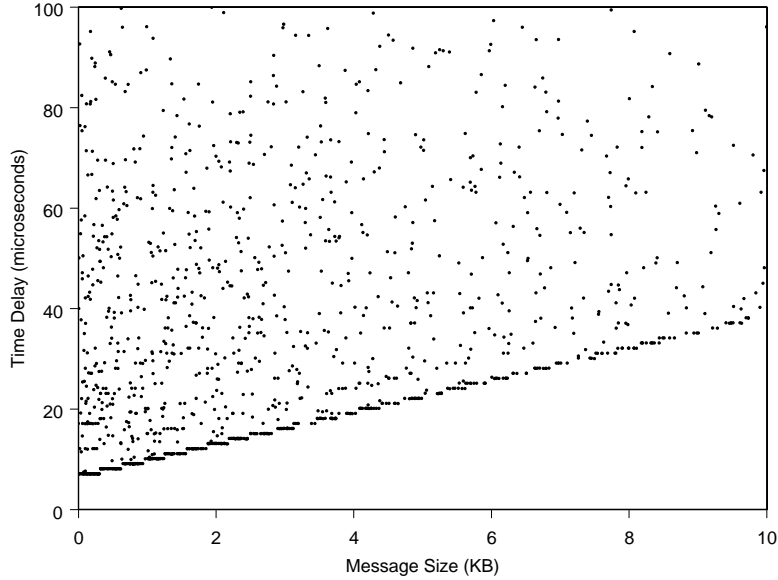


Figure 14: Example output data from a 64-node simulation where each node generates an average of 100 MB/s of message traffic on the network. The mean message size is 4 KB.

the workload is highly generalized. Nevertheless, the simulation does show typical characteristics for a moderately-loaded system.

We have also used the simulation to compare two typical types of application communication patterns. Our “uniform” pattern represents applications that send messages at random between computational nodes: specifically, each computational node has equal probabilities of sending a message to any other computational node. Figure 15 illustrates the sort of network traffic such a workload creates in our simulation. Our “CFD” pattern represents applications that send messages to neighbors on a computational grid, such as a computational fluid dynamics application does: specifically, the computational nodes are assigned locations on a three dimensional cube (e.g., on a $4 \times 4 \times 4$ cube for a 64-node simulation or on a $16 \times 16 \times 16$ cube for a 4096-node simulation), and each node has equal probabilities of sending a message only to its immediate neighbors on the grid. Figure 16 illustrates that sort of network traffic in our simulation. The CFD pattern clearly uses the higher layers of the fat tree much less than the uniform pattern.

6 Medium-Fidelity Quadrics Network Simulation Design and Implementation

The *à la carte* medium-fidelity network model builds on the low-fidelity model discussed in the previous section by enhancing its accuracy and realism. This model is much closer to representing actual hardware and mimicking the behavior of network protocols in use on real systems. We expect the resolution of this representation to be sufficient for even the most detailed network studies.

6.1 Requirements

Our primary requirement is the ability to accurately model the movement of packets in Quadrics networks consisting of Elan network interface cards [14, 15, 16] connected to Elite switches [17] in a fat-tree network at nearly *flit* (16-bit units) resolution. Figure 17 illustrates the sequence

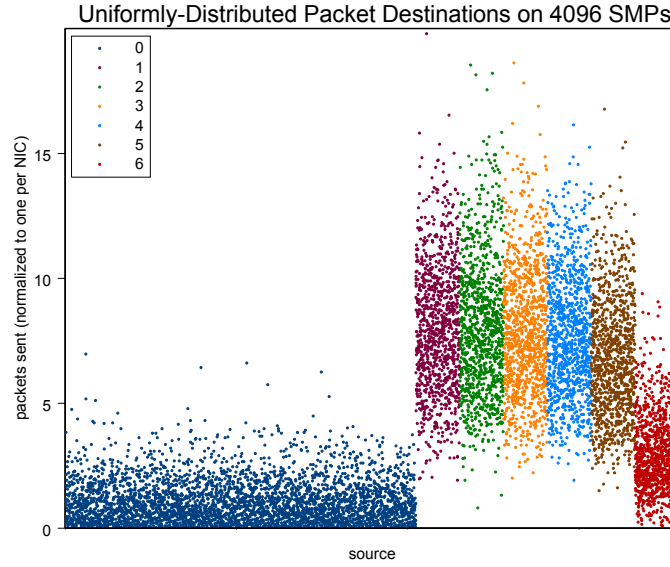


Figure 15: Number of message packets for the “uniform” pattern of workload communication in a 4096-node simulation: The horizontal axis represents the simulation entity (nodes in dark blue, and switches in other colors for each level) and the vertical axis represents the number of packets sent by the entity, normalized on a scale so that the nodes send an average of one packet.

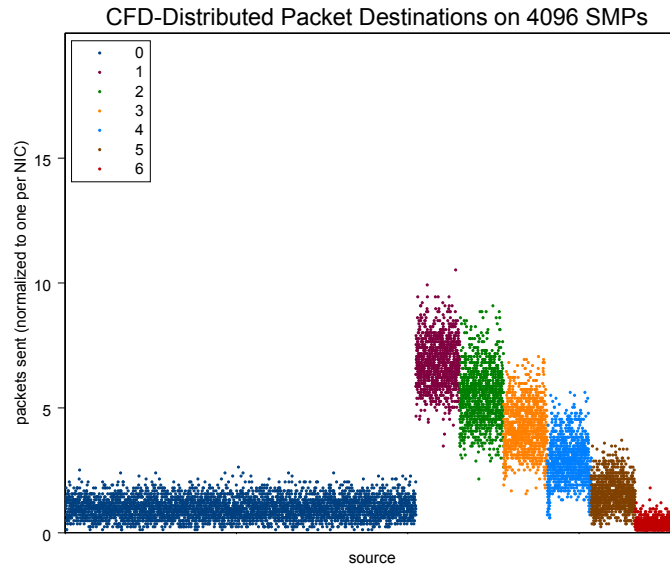


Figure 16: Number of message packets for the “CFD” pattern of workload communication in a 4096-node simulation: The horizontal axis represents the simulation entity (nodes in dark blue, and switches in other colors for each level) and the vertical axis represents the number of packets sent by the entity, normalized on a scale so that the nodes send an average of one packet.

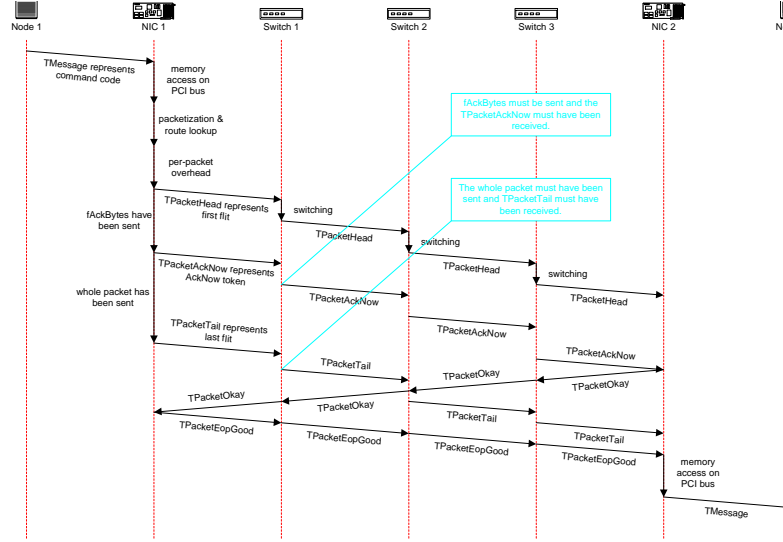


Figure 17: Steps in the movement of a message from one computational node to another in a Quadrics network of Elan network interface cards and Elite switches.

of operations it takes to move a message from one computational node to another through this sort of network. We need to accurately track the movement of the message across the PCI bus between main memory and the NIC, account for its packetization, and clock the transfer of data across the network. Because contention may exist in the network, different parts of the packet may move at different speeds through the switches (i.e., buffering and delays may occur anywhere in the network).

In addition to more accurate representation of network, we also need additional tools to assemble networks that are not perfect fat-trees. This includes linking several fat trees together, and handling trees that are thinned (i.e., have fewer switches) in their upper layers. It is also necessary to have the ability to input cable lengths, as these range over orders of magnitude (from centimeters to meters) in real networks.

6.2 Design

It is a fairly simple matter to augment the design for our low-fidelity network representation to the higher fidelity required. The basic **TSMPNode**, **TNIC**, and **TSwitch** entities remain essentially unchanged. We use object factory methods to generate the appropriate **TRoutingAlgorithm** or **TFlowControl** instances for the fidelity specified in the DML input file. The medium fidelity **TCircuitAlgorithm** and **TCircuitControl** classes are similar to their low fidelity counterparts except for their internal logic, which is considerably more complex, and the type of packets they handle. The design relies on the tracking of the head and tail of the packet throughout its history, along with various flit-level tokens specified in the Elan protocol. Figure 18 shows the class diagram for the medium-fidelity simulation.

In addition to tracking the movement of the head and the tail of packets through the NICs and switches, we account for the existence of two virtual channels sharing bandwidth at switches (but without age-based priorities, etc.). The Elan *AckNow* request may occur anywhere within the packet (usually after 64 bytes or at the end of the packet). The *EOP_GOOD* tokens free the virtual channels used by the packet. The *STARTx/STOPx* tokens are accounted for by buffering of packets

at the incoming links to switches if no outgoing virtual channel is available. We model the PCI bus as half-duplex, and account for writes to the NIC’s command port. Finally, we allow wildcards for packet routing on upward links. We do not yet model error conditions. We believe this is nearly the highest fidelity we can achieve without simulating at the flit level, which would be prohibitively slow. (Adding further resolution may not be cost effective because of the uncertainties involved in the performance of the operating system on the node, memory issues, and PCI bus behavior.)

The basic strategy for dealing with packets is as follows: When the head of the packet reaches an entity like a NIC, switch, or node, it leaves a **TPacketReservation** entity at the entity. The head of the packet is forwarded along the route as soon as possible—it might be delayed slightly for switch logic or may be delayed significantly if it is queued for later transmission. As soon as the head leaves the entity, the reservation starts keeping track of how many bytes remain to be transmitted. The tail of the packet cannot be forwarded until it has been received from the previous stage and the number of bytes remaining is zero. The “okay” event proceeds along the reverse path at full speed, and the “good” event cleans up the reservations. There is some fairly complex timekeeping logic for multiplexing the transmission of packets in switches and the receipt of them in NICs.

6.3 Implementation

We track the leading edge (“head”), trailing edge (“tail”), *AckNow* token, *PACK_OK* token, and *EOP_GOOD* token for packets in the network, which are implemented as **TPacketHead**, **TPacketTail**, **TPacketAckNow**, **TPacketOkay**, and **TPacketEopGood** subclasses of our existing **TPacket** class (see Fig. 18). The **TPacket** class differs from its low-fidelity counterpart in that it supports virtual channels and is tailored for its subclasses that track flit- level events. In Figure 19 we outline the following scenario for the transmission of a message from one computational node to another.

1. A **TWorkloadSender** process creates a **TMessage** event and writes it to the **TSMPNode.fOutBus** channel after a delay of **TSMPNode.fCommunicationsLatency** seconds. This represents the command code being written from main memory. The **TSMPNode.fCommunicationsLatency** parameter comes from the DML and represents the amount of time the communications libraries (MPI, TPORTS, ELANLIB, and ELAN3LIB) take to handle the message.
2. The **TMessage** event moves across the link from **TSMPNode.fOutBus** to **TNIC.fInBus**, acquiring the delay specified in the DML. This represents the command code being written to the NIC’s memory. The delay must equal the reciprocal of the PCI bus bandwidth.
3. The **TWorkloadListener** process queues the **TMessage**. This represents the request for a DMA being queued.
4. The **TNICSender** process pops the top **TMessage** off its queue of message requests. It waits to obtain a lock on the **TNIC.fBus** semaphore. The semaphore accounts for the fact that the PCI bus is half-duplex. Obtaining the lock represents the first byte of the DMA transfer.
5. Once the **TNICSender** has the lock, it delays for **TMessage.fSize / TNIC.fBusBandwidth** seconds and then releases the lock. This write represents the last byte of the DMA transfer. The **TNIC.fBusBandwidth** parameter represents the bandwidth of the PCI bus.
6. The **TNICSender** packetizes the **TMessage** by creating a queue of **TPacketReservation** instances. When these **TPacketReservation** instances are created each holds a **TPacketHead**, **TPacketTail**, and **TPacketAckNow** event and has a **TPacketReservation.fBytesRemaining** attribute equal to the packet size. The **TPacketHead.fAckBytes** attribute should be set according to when *PACK_OK* should

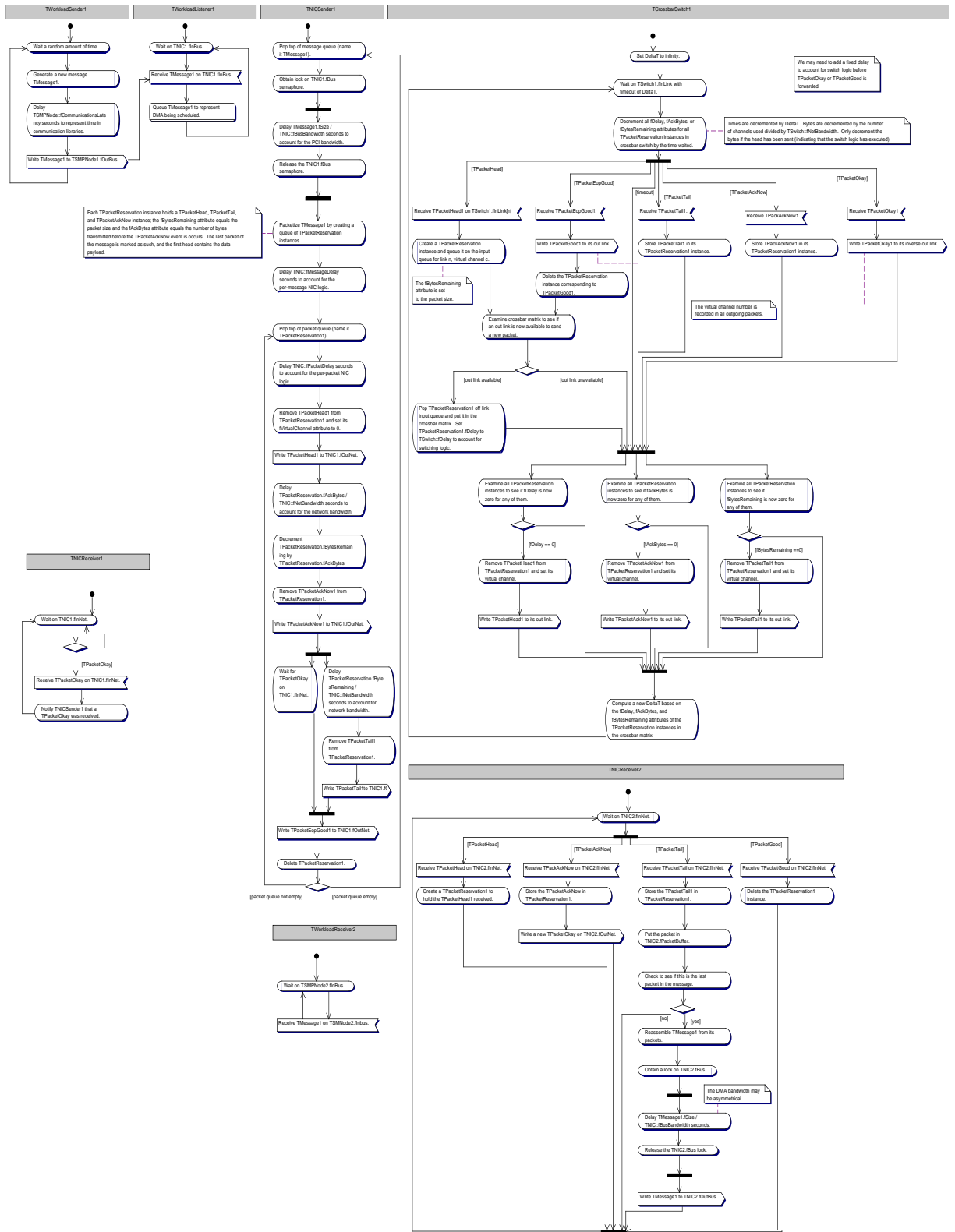


Figure 19: UML activity diagram outlining the processing of messages and packets in the medium-fidelity Quadrics network simulation.

be sent by the receiver: this is the number of bytes that must be sent before **TPacketAckNow** is sent. We delay **TNIC.fMessageDelay** seconds here to account for the NIC logic that executes.

7. The **TNICSender** starts transmitting the packet at the top of its queue by writing the **TPacketHead** event to the **TNIC.fOutNet** channel after a delay of **TNIC.fPacketDelay** seconds. Before being written, **TPacketHead.fVirtualChannel** is set to zero. Since the NIC can only send on virtual channel 0, it does not have to share bandwidth between both virtual channels.
8. The **TNICSender** waits **TPacketReservation.fAckNowBytes / TNIC.fNetBandwidth** seconds and then writes the **TPacketAckNow** event to the **TNIC.fOutNet** channel. The **TNIC.fNetBandwidth** parameter represents the bandwidth of the network links.
9. The **TNICSender** waits **TPacketReservation.fBytesRemaining / TNIC.fNetBandwidth** seconds and then writes the **TPacketTail** event to the **TNIC.fOutNet** channel.
10. After **TPacketTail** has been sent and the **TPacketOkay** event from the **TNICReceiver** has been received, the **TNICSender** sends a **TPacketEopGood** event.
11. Meanwhile, the **TPacket...** events move from the **TNIC.fOutNet** channel to the **TSwitch.fInNet**, acquiring a delay specified in the DML. The delay corresponds to the physical speed of transmission in the cable.
12. When a **TPacketHead** event is received at a **TSwitch**, a **TPacketReservation** instance is created and queued on the input buffer corresponding to its input port and virtual channel. The **TPacketReservation** instance holds the **TPacketHead** and its **fBytesRemaining** field is set to **TPacketHead.fSize**. Each switch has 8×2 input buffers. This buffering represents the *STOP_x* and *START_x* of the flit-level protocol.
13. Whenever a **TSwitch** has queued a **TPacketHead** or handled a **TPacketEopGood** event (described below), it next checks if the **TPacketReservation** at the top of the corresponding input buffer can be served by seeing whether one of the virtual channels for its output port is free. The 16×16 crossbar matrix holds pointers to **TPacketReservation** instances currently being transmitted: each row (input port/channel) and column (output port/channel) may have at most one non-null entry. If a virtual channel is available for the output port, then a pointer to the **TPacketReservation** is set there and a delay of **TSwitch.fDelay** seconds is assigned to **TPacketReservation.fDelay** to account for switch logic. Note that a wildcard in the packet's route can be assigned to any upward output port/channel.
14. A **TSwitch** moves time forwards as follows: Find the minimum of the **TPacketReservation.fDelay** values for transmitting packets with this not zero and **TPacketReservation.fBytesRemaining * channelsUsed[outport] / TSwitch.fNetBandwidth** values for transmitting packets with **TPacketReservation.fDelay** zero, where **channelsUsed** is the number of virtual channels being used on a particular output port. Call **SSF.Entity.waitFor(SSF.InChannel** channels, ltime.t timeout)** with that minimum value as the timeout parameter. Subtract the actual time waited from the **TPacketReservation.fDelay** attributes or **channelsUsed[outport] / TSwitch.fNetBandwidth** from the **TPacketReservations.fBytesRemaining** attribute, as appropriate. At this stage either a new event has arrived on an input port, or one of the **TPacketReservation.fDelay** or **TPacketReservation.fBytesRemaining** attributes is zero.
15. If **TPacketReservation.fDelay** is zero, then remove the **TPacketHead** and write it to the proper **TSwitch.fOutLink** channel. Before being written, **TPacketHead.fVirtualChannel** is set to the correct number and the **channelsUsed[outport]** is incremented. Go to step 14.

16. If **TPacketReservation.fAckNowBytes** is zero, then remove the **TPacketAckNow** if it exists and write it to the proper **TSwitch.fOutLink** channel. Go to step 14.
17. If **TPacketReservation.fBytesRemaining** is zero, then remove the **TPacketTail** if it exists and write it to the proper **TSwitch.fOutLink** channel. Go to step 14.
18. If a **TPacketAckNow** event has been received, hold it in the corresponding **TPacketReservations** instance in the **TSwitch**. If the **TPacketReservation.fAckNowBytes** attribute is zero, then go to step 16; otherwise, go to step 14.
19. If a **TPacketTail** event has been received, hold it in the corresponding **TPacketReservations** instance in the **TSwitch**. If the **TPacketReservation.fBytesRemaining** attribute is zero, then go to step 17; otherwise, go to step 14.
20. If a **TPacketEopGood** event has been received, then delete the **TPacketReservation** instance and write the **TPacketEopGood** event to the proper **TSwitch.fOutLink** channel. Go to step 14.
21. If a **TPacketOkay** event has been received, then write it to the proper **TSwitch.fOutLink** channel of the inverse crossbar matrix. Go to step 14.
22. When the **TNICReceiver** process receives a **TPacketHead** event, it creates a **TPacketReservation** instance that holds the **TPacketHead**.
23. If a **TPacketAckNow** event has arrived, then hold this in the **TPacketReservation** instance and send a **TPacketOkay** event on the virtual channel **TPacketHead.fVirtualChannel**.
24. If a **TPacketTail** event has arrived, then hold this in the **TPacketReservation** instance and add the packet to **TNICReceiver.fPktBuffer**. If this was the last packet in the message, go to step 26.
25. If a **TPacketEopGood** event has arrived, then delete the **TPacketReservation** instance.
26. Once the **TNICReceiver** has all of the packets, it reconstitutes it into a message and waits to obtain a lock on the **TNIC.fBus** semaphore. Once the **TNICReceiver** has the lock, it delays for **TMessage.fSize / TNIC.fBusBandwidth** seconds and then releases the lock. The **TMessage** is written to **TNIC.fOutBus**.
27. The **TWorkloadReceiver** process handles the receipt of the **TMessage**.

Our implementation of more advanced simulation output data collection features relies on a new **TDataCollector** class that mediates all types of data collection and error logging. In addition to the complete event detail available in the low-fidelity model—which tracks the history of messages and propagation of packets—it can support the collection of summary data which may be sliced into time windows. Summary data includes information on queue sizes, throughputs at switches, port usages, timeouts, and path lengths. The output is configurable so that it can be retrieved either as it is generated or after the simulation complete.

Addition of a ping capability to our simple workload is accomplished by adding another DML parameter to the “node” section to specify whether the messages have a fixed size and interval or are chosen from a distribution. This allows us to perform detailed validation studies for the simulation.

7 Statistical Characterization of Workloads

The theoretical and computational issues involved in large-scale network modeling have direct analogues in some of the most challenging and important problems in statistical physics. The complicating features which appear in both contexts include multiple spatial and temporal scales, and strongly correlated dynamics and stochasticity. As a result we are bringing to bear on problems arising in network performance modeling techniques originally developed and already highly successful in the context of nonequilibrium statistical physics.

One such technique is a numerical Rayleigh-Ritz method [18]. This is a variational formulation for the time dependence of the probability distribution functions of network variables. Whether one uses a discrete event dynamical system or stochastic fluid level approach, the network variables are described by a large system of nonlinear, coupled, stochastic equations. Due to their nonlinear nature, these equations cannot be solved exactly and require that approximations be made. The approximations which describe higher order statistics in terms of lower ones are known as a closure. The Rayleigh-Ritz variational method then tests these statistical closure ideas using the exact dynamics (i.e., the original equations), but more cheaply and quite often under more extreme circumstances than is feasible by direct numerical simulation.

Closure schemes can then be developed by using educated guesses for the system statistics. These “guesses” may be inspired from an analysis of the data using the visualization tool. As a result, empirical data from actual networks and workloads may be exploited in a numerical calculation. Additional statistical information can be calculated within the variational formulation, including multi-time correlations, and the probability of large fluctuations in performance. Moreover, there are internal consistency checks available which may be used as diagnostics to detect *a priori* faulty predictions of the closures, potentially reducing the amount of time spent on inadequate models.

Using this methodology we are attempting to address several types of problems such as the probability of buffer overflow and likelihood of long-time performance averages. Moreover, we expect that well chosen closures will lead to being able to answer these questions with greatly reduced computational effort.

7.1 Preliminary Results

We have explored our preliminary data sets (see the low fidelity simulation results on page 27) for interesting statistical structure using the method of principal components analysis (PCA). We postprocess our simulation output data to construct a vector $\mathbf{q}(t)$ with components $q_i(t)$ equal to the number of messages waiting to be transmitted through the network from computational node $i = 1 \dots n$ at time step $t = 0 \dots T$. We call $\mathbf{q}(t)$ the *queue length* time series. From this we compute a covariance matrix

$$\mathbf{C} = \frac{1}{T} \sum_{t=0}^T [\mathbf{q}(t) - \bar{\mathbf{q}}] [\mathbf{q}(t) - \bar{\mathbf{q}}]^t, \quad (11)$$

where $\bar{\mathbf{q}} = \frac{1}{T+1} \sum_{t=0}^T \mathbf{q}(t)$, and find its eigenvalues λ_i ordered such that $\lambda_i \geq \lambda_{i+1}$. In general, some of the eigenvalues can be attributed to random noise rather than to the correlation structure in the queue length time series. To determine the number of statistically significant eigenvalues, k , we consider random permutations of the original series:

$$q'_i(t) = q_i(t')|_{t'=\mathbb{P}t} \quad (12)$$

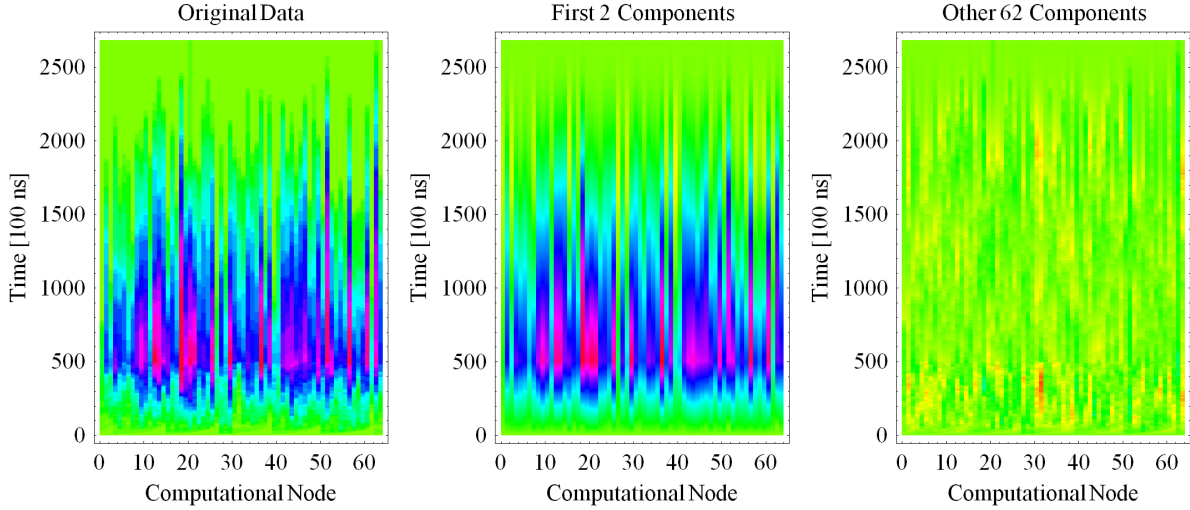


Figure 20: Queue lengths in a simulation with 64 computational nodes (left) and the contribution from the two eigenmodes with statistically significant eigenvalues (middle) and the remaining 62 eigenmodes (right): The horizontal axis represents node number, the vertical axis is simulated time, and colors range from green (queue length of zero) to red (queue length of twelve).

where \mathbb{P} is an operator that maps $\{0, \dots, T\}$ to a random permutation of $\{0, \dots, T\}$. We can find the largest eigenvalue λ'_1 for several realizations of the series $\mathbf{q}'(t)$ to determine the cut-off below which eigenvalues λ_i can be attributed to random noise in the correlation matrix.

Preliminary results indicate that only a few eigenmodes, k , are statistically significant for the “CFD” and “uniform” workloads (see page 27) if the simulated network is operated near its capacity. Figure 20 shows how $k = 2$ modes are sufficient to explain most of the variation in queue length for a simulation with $n = 64$ computational nodes. We hope that further analysis using the workloads generated by real applications on our medium-fidelity network will show similarly encouraging results.

8 Conclusion

We have outlined our low- and medium-fidelity fat-tree network models along with the simulation and visualization technologies supporting them. Our ongoing work focuses on validating these models and developing more faithful representations of workload to be used in conjunction with them. Some progress has already been made in the statistical modeling of workloads. Our component-based development process will enable us to seamlessly compose hardware, protocols, and workloads of varying fidelities into a single simulation.

The major risk in our effort is that of scaling: Simulating extreme scale systems is resource-intensive even with efficient simulator fidelity. Constructing high-fidelity models of processors or networks is a labor-intensive process and modeling operating system behavior might be required in some studies. Other risk areas are portability (accurately measuring directly-executed applications requires non-portable timers, and the modeling of low-level networking APIs may reduce portability) and validation (detailed measurements of applications on large machines are needed to validate a simulation).

9 Acknowledgements

This work was carried out under the auspices of the Department of Energy at Los Alamos National Laboratory under ASCI DisCom2. We would like to thank LANL's DisCom2 project leader, Stephen Turpin, for his support.

Thanks to members of the extended LANL ParSim team for technical input, including Mike Boorman, Fabrizio Petrini, and Harvey Wasserman; and to David Nicol and Jason Liu (Dartmouth College), and Josep Torrellas (University of Illinois at Urbana-Champaign). Anand Sivasubramaniam (Pennsylvania State University) and Madhav Marathe participated in the initial project proposal.

References

- [1] A. Hoisie, O. Lubeck, H. Wasserman, "Scalability Analysis of Multidimensional Wavefront Algorithms on Large-Scale SMP Clusters," *Proceedings of the 7th Symposium on the Frontiers of Massively Parallel Computations, Frontiers* February (1999).
- [2] A. Hoisie, O. Lubeck, H. Wasserman, "Performance Analysis of Wavefront Algorithms on Very-Large Scale Distributed Systems," *Lecture Notes in Control and Information Sciences* 249, Springer-Verlag, (1999).
- [3] D. E. Culler, R. Karp, D. Patterson, A. Sahay, K. E. Schauser, E. Santos, R. Subramonian, and T. von Eiken, *LogP: Towards a Realistic Model of Parallel Computation*. Proceedings of the Fourth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, San Diego, CA, May 1993.
- [4] M. Goudreau, K. Lang, S. Rao, T. Suel, and T. Tsantilas, *Towards Efficiency and Portability: Programming with the BSP Model* Proceedings of the 8th Annual ACM Symposium on Parallel Algorithms and Architectures, Padova, Italy, 1996.
- [5] R. Kaufman, *The Q Supercomputer and Compaq*. High Performance Technical Computing News, Issue 18, November 2000, http://www.compaq.com/hpc/news/news_hpc_60171.html.
- [6] Richard M. Fujimoto. *Parallel and Distributed Simulation Systems*, John Wiley & Sons, Inc., 2000.
- [7] <http://www.ssfnet.org>.
- [8] James H. Cowie, David M. Nicol, Andy T. Ogielski. *Modeling the Global Internet*. Computing in Science & Engineering 1(1):30-38, 1999.
- [9] Jason Liu and David M. Nicol. *Dartmouth Scalable Simulation Framework User's Manual*. Dartmouth College Dept. of Computer Science, February 6, 2002.
- [10] S. Das, R. Fujimoto, K. Panesar, D. Allison, and M. Hybinette. Gtw: A time warp system for shared memory multiprocessors. In *Winter Simulation Conference*, 1994.
- [11] R. Bagrodia, R. Meyer, M. Takai, Y. Chen, X. Zheng, J. Martin, and H. Song. Parsec: A Parallel simulation environment for complex systems. *IEEE Computer*, 1998.
- [12] S. Prakash and R. Bagrodia. MPI-SIM: Using parallel simulation to evaluate MPI programs. In *Winter Simulation Conference*, 1998.

- [13] J. Liu, D. Nicol, B. Premore, A. Poplawski, "Performance Prediction of a Parallel Simulator", Proc. Parallel and Distributed Simulation Conf. Atlanta, 1999.
- [14] "Elan Kernel Communications Manual"
- [15] "Elan Programming Manual"
- [16] "Elan Reference Manual"
- [17] "Elite Reference Manual"
- [18] F. J. Alexander and G. L. Eyink. *A Rayleigh-Ritz Calculation of the Effective Potential far from Equilibrium*. Physical Review Letters 78 1, 1997.
- [19] David M. Nicol and Jason Liu. *Composite Synchronization in Parallel Discrete-Event Simulation*. IEEE Transactions on Parallel and Distributed Systems 13(5), May, 2002. To appear.
- [20] Y.-H. Low, C.-C. Lim, W. Cai, S.-Y. Huang, W.-J. Hsu, S. Jain, and S. Turner. *Survey of Languages and Runtime Libraries for Parallel Discrete-Event Simulation* Simulation, 72(3):170-186, March, 1999.
- [21] <http://www.cs.dartmouth.edu/research/DaSSF>
- [22] <http://www.cc.gatech.edu/computing/pads/tech-parallel-gtw.html>
- [23] <http://www.cc.gatech.edu/computing/compass/pdns>
- [24] <http://www.isi.edu/nsnam/ns>
- [25] L. Breslau, D. Estrin, K. Fall, S. Floyd, J. Heidemann, A. Helmy, P. Huang, S. McCanne, K. Varadhan, Y. Xu, H. Yu. *Advances in Network Simulation* IEEE Computer, 33(5):59-67, May, 2000.
- [26] S. Mukherjee, S. Reinhardt, B. Falsafi, M. Litzkow, M. Hill, D. Wood, S. Huss-Lederman, and J. Larus. *Wisconsin Wind Tunnel II: A Fast, Portable Parallel Architecture Simulator* IEEE Concurrency, Oct-Dec, 2000.
- [27] R. L. Bagrodia. *Parallel Languages for Discrete-Event Simulation Models* IEEE Computational Science & Engineering, Apr-June, 1998.
- [28] R. Bagrodia, E. Deelman, and T. Phan. *Parallel simulation of Large-Scale Parallel Applications* Intl. J. of High Performance Computing Applications, 15(1):3-12, Spring, 2001.
- [29] A. Srivastava and A. Eustace. *ATOM: A System for Building Customized Program Analysis Tools* WRL Research Report 94/2, March, 1994.
- [30] <http://icl.cs.utk.edu/projects/papi>
- [31] <http://www.fz-juelich.de/zam/PCL>
- [32] <http://www.ittc.ku.edu/utime>
- [33] <http://www.ahpcc.unm.edu/homunculus/indexold.html>